

1. Utilizar el siguiente set de datos para calcular paso por paso.

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: $\bar{X} = \frac{\sum_{i=1}^n X_i}{n}$ } $n=21$

• Para X_1

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

$$\Rightarrow \bar{X}_1 = \frac{49}{21} \Rightarrow \therefore \bar{X}_1 = 2,333 //$$

• Para X_2

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

$$\Rightarrow \bar{X}_2 = \frac{93}{21} \Rightarrow \therefore \bar{X}_2 = 4,428 //$$

• Para X_3

$$\rightarrow \bar{X}_3 = (28+24+30+32+18+41+49+5+18+1+62+17+24+19+42+56+44+64+17+17) \cdot \frac{1}{21}$$

$$\Rightarrow \bar{X}_3 = \frac{659}{21} \Rightarrow \therefore \bar{X}_3 = 31,38 //$$

Mediana: Valor central

• Para X_1 :

→ $Me_1 = 0, 0, 1, 1, 1, 1, 1, 1, 2, 2, \textcircled{2}, 3, 3, 3, 3, 3, 4, 4, 4, 5, 5$

$$\therefore Me_1 = 2 //$$

• Para X_2 :

→ $Me_2 = 0, 1, 2, 2, 2, 3, 3, 3, 3, 3, \textcircled{4}, 5, 6, 6, 6, 6, 8, 8, 9, 9$

$$\therefore Me_2 = 4 //$$

• Para X_3 :

→ $Me_3 = 1, 5, 17, 17, 17, 18, 18, 19, 24, 24, \textcircled{28}, 30, 32, 41, 42, 44, 44, 56, 56, 62, 64$

$$\therefore Me_3 = 28 //$$

Desviación Estándar $\sigma = \sqrt{\frac{1}{N} \cdot \sum_{i=1}^N (x_i - \bar{x})^2}$

$\text{Var}(x)$

• Para X_1

$$\Rightarrow \text{Var}(X_1) = [(4-2,333)^2 + (2-2,333)^2 + (2-2,333)^2 + (3-2,333)^2 + (1-2,333)^2 + (3-2,333)^2 + (3-2,333)^2 + (0-2,333)^2 + (1-2,333)^2 + (0-2,333)^2 + (5-2,333)^2 + (1-2,333)^2 + (2-2,333)^2 + (1-2,333)^2 + (3-2,333)^2 + (4-2,333)^2 + (4-2,333)^2 + (3-2,333)^2 + (5-2,333)^2 + (1-2,333)^2 + (1-2,333)^2] \cdot \frac{1}{21}$$

$$\Rightarrow \text{Var}(X_1) = 2,222 \Rightarrow \sigma = \sqrt{2,222} \Rightarrow \therefore \sigma = 1,49 //$$

• Para X_2

$$\Rightarrow \text{Var}(X_2) = [(4-4,428)^2 + (2-4,428)^2 + (4-4,428)^2 + (5-4,428)^2 + (3-4,428)^2 + (6-4,428)^2 + (6-4,428)^2 + (1-4,428)^2 + (3-4,428)^2 + (0-4,428)^2 + (9-4,428)^2 + (2-4,428)^2 + (3-4,428)^2 + (6-4,428)^2 + (8-4,428)^2 + (8-4,428)^2 + (6-4,428)^2 + (9-4,428)^2 + (2-4,428)^2 + (2-4,428)^2] \cdot \frac{1}{21}$$

$$\Rightarrow \text{Var}(X_2) = 6,53 \Rightarrow \sigma = \sqrt{6,53} \Rightarrow \therefore \sigma = 2,555 //$$

• Para X_3

$$\begin{aligned} \text{Var}(X_3) = & [(28-31,35)^2 + (24-31,35)^2 + (30-31,35)^2 + (32-31,35)^2 \\ & + (16-32,35)^2 + (41-32,35)^2 + (44-32,35)^2 + (5-32,35)^2 + (18-32,35)^2 \\ & + (1-32,35)^2 + (62-32,35)^2 + (17-32,35)^2 + (24-32,35)^2 + (19-32,35)^2 + (42-32,35)^2 \\ & + (56-32,35)^2 + (56-32,35)^2 + (44-32,35)^2 + (64-32,35)^2 + (17-32,35)^2 \\ & + (17-32,35)^2] \cdot \frac{1}{21} \end{aligned}$$

$$\rightarrow \text{Var}(X_3) = 314,807 \Rightarrow \sigma = \sqrt{314,807} \Rightarrow \therefore \sigma = 17,742 //$$

Moda: El valor que mas se repite

Para X_1

$$Mo_1 = 1 //$$

↓
6 repeticiones

Para X_2

$$Mo_2 = 3 //$$

↓
5 repeticiones

Para X_3

$$Mo_3 = 17 //$$

↓
3 repeticiones

1.2. Dibujar un boxplot a mano

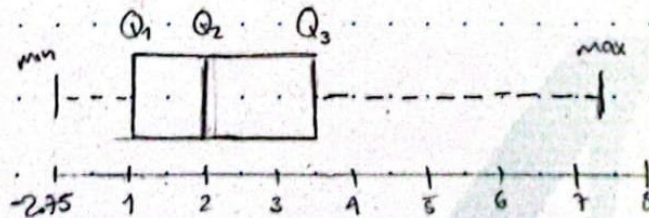
2. Boxplot

- $IQR = Q_3 - Q_1$
- Q_1 = Mediana de la mitad inferior de los datos
- Q_3 = " " " Superior " "
- Barrera Superior = $Q_3 + 1.5 IQR$
- Barrera Inferior = $Q_1 - 1.5 IQR$
- Q_2 = El 50% de los datos (~~promedio~~) Mediana

$$Q_1 = \frac{1+1}{2} \Rightarrow 1 \quad Q_2 = 2 \quad Q_3 = \frac{3+4}{2} \Rightarrow 3.5 \quad IQR = 3.5 - 1 \Rightarrow 2.5$$

$$\text{Barrera Inferior} = 1 - 1.5 \times 2.5 \Rightarrow -2.75$$

$$\text{Barrera Superior} = 7.25$$



Para X_2 :

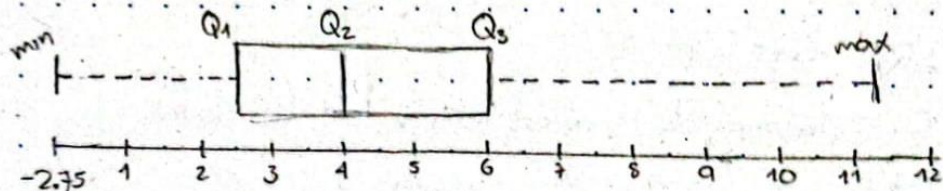
$$Q_1 = \frac{3+2}{2} \Rightarrow 2.5 \quad Q_3 = 6 \quad Q_2 = 4 \quad IQR = 6 - 2.5 = 3.5$$

Barrera Inferior

Barrera Superior

$$2.5 - 1.5 \times 3.5 \Rightarrow -2.75$$

$$6 + 1.5 \times 3.5 \Rightarrow 11.25$$



Para X_3 :

$$Q_1 = \frac{17+18}{2} \Rightarrow 17.5 \quad Q_3 = 44 \quad Q_2 = 28 \quad IQR = 44 - 17.5 = 26.5$$

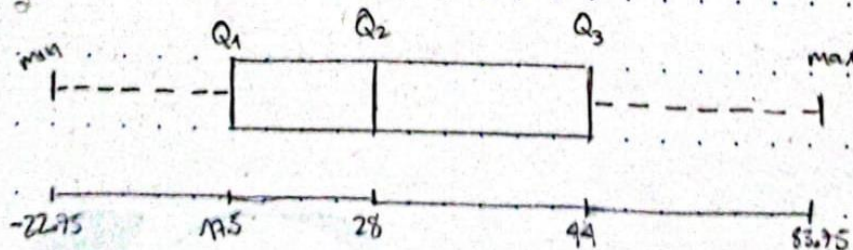
Barrera Inferior

Barrera Superior

$$17.5 - 1.5 \times 26.5 \Rightarrow -22.75$$

$$44 + 1.5 \times 26.5 \Rightarrow 83.75$$

Lo 2



1.3. Cuál 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
$$\text{Cov}(X_1, X_2) = \frac{\sum (X_{1i} - \bar{X}_1) \cdot (X_{2i} - \bar{X}_2)}{N}$$

$$\Rightarrow \text{Cov}(X_1, X_2) = \left[(4-2,333)(4-4,428) + (2-2,333)(3-4,428) + (2-2,333)(4-4,428) \right. \\
 + (3-2,333)(5-4,428) + (1-2,333)(3-4,428) + (3-2,333)(6-4,428) \\
 + (3-2,333)(6-4,428) + (0-2,333)(1-4,428) + (1-2,333)(3-4,428) \\
 + (0-2,333)(0-4,428) + (5-2,333)(9-4,428) + (1-2,333)(2-4,428) \\
 + (2-2,333)(3-4,428) + (1-2,333)(3-4,428) + (3-2,333)(6-4,428) \\
 + (4-2,333)(8-4,428) + (4-2,333)(8-4,428) + (3-2,333)(6-4,428) \\
 \left. + (5-2,333)(9-4,428) + (1-2,333)(2-4,428) + (1-2,333)(2-4,428) \right] \cdot \frac{1}{21}$$

$$\Rightarrow \text{Cov}(X_1, X_2) = \frac{75}{21} \Rightarrow \therefore \text{Cov}(X_1, X_2) = 3,571 //$$

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

$$\Rightarrow \text{Cor}(X_1, X_2) = \frac{75}{\sqrt{46,662} \cdot \sqrt{137,13}} \Rightarrow \therefore \text{Cor}(X_1, X_2) = 0,937 //$$

1.5. Explica la relación entre covarianza y correlación.

La covarianza explica que tipo de comportamiento tienen dos variables en conjunto, es decir, si el resultado es positivo quiere decir que si una aumenta la otra tiende a aumentar; si el resultado es negativo quiere decir que si una aumenta la otra tiende a disminuir.

Como se puede observar, su descripción concuerda por mucho con la correlación ya que también analiza como es el comportamiento de una variable respecto a la otra, en otras palabras, si existe una dependencia de una variable hacia la otra.

En lo único que difiere es en la interpretación, ya que la covarianza no esta normalizada, lo que puede dificultar la interpretación de esta y comparación entre las diferentes escalas de las variables, cosa que no pasa en la correlación, ya que esta última si esta normalizada (porque se está dividiendo básicamente por sus respectivas varianzas), es decir que sus datos varían entre -1 y 1, lo que ayuda significativamente a interpretar los valores.

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 $K=3$

Inicio

- Grupo 3 avg = -
 - o Centroide $X_1 = \frac{4+2+2+3+3+1+4}{7} \Rightarrow 3,142$
 - o Centroide $X_2 = \frac{4+3+4+5+6+8+8}{7} \Rightarrow 5,428$
- Grupo 2 avg = -
 - o Centroide $X_1 = \frac{1+3+3+0+1+5+1+1}{8} \Rightarrow 1,875$
 - o Centroide $X_2 = \frac{3+6+6+1+3+9+2+2}{8} \Rightarrow 4$
- Grupo 1 avg = -
 - o Centroide $X_1 = \frac{0+5+1+2+1+3}{6} \Rightarrow 2$
 - o Centroide $X_2 = \frac{0+9+2+3+3+6}{6} \Rightarrow 3,833$

$dis_{-}C_3 = 1 \cdot \sqrt{(4-3,142)^2 + (4-5,428)^2} \Rightarrow 1,665$ $5 \cdot \sqrt{(1-3,142)^2 + (1-5,428)^2} = 3,238$
 $2 \cdot \sqrt{(2-3,142)^2 + (2-5,428)^2} \Rightarrow 2,684$ $6 \cdot \sqrt{(3-3,142)^2 + (6-5,428)^2} = 0,589$
 $3 \cdot \sqrt{(2-3,142)^2 + (4-5,428)^2} \Rightarrow 1,629$ $7 \cdot \sqrt{(3-3,142)^2 + (6-5,428)^2} = 0,589$
 $4 \cdot \sqrt{(3-3,142)^2 + (5-5,428)^2} \Rightarrow 0,4517$ $8 \cdot \sqrt{(0-3,142)^2 + (1-5,428)^2} = 5,130$
 $5 \cdot \sqrt{(3-3,142)^2 + (6-5,428)^2} \Rightarrow 0,589$ $9 \cdot \sqrt{(1-3,142)^2 + (3-5,428)^2} = 3,238$
 $6 \cdot \sqrt{(4-3,142)^2 + (8-5,428)^2} \Rightarrow 2,710$ $10 \cdot \sqrt{(0-3,142)^2 + (0-5,428)^2} = 6,272$
 $7 \cdot \sqrt{(4-3,142)^2 + (8-5,428)^2} \Rightarrow 2,710$ $11 \cdot \sqrt{(5-3,142)^2 + (9-5,428)^2} = 4,025$
 $8 \cdot \sqrt{(1-3,142)^2 + (3-5,428)^2} \Rightarrow 1,613$ $12 \cdot \sqrt{(3-3,142)^2 + (6-5,428)^2} \Rightarrow 0,589$
 $9 \cdot \sqrt{(2-3,142)^2 + (3-5,428)^2} \Rightarrow 2,684$ $13 \cdot \sqrt{(5-3,142)^2 + (9-5,428)^2} \Rightarrow 4,025$
 $10 \cdot \sqrt{(1-3,142)^2 + (1-5,428)^2} \Rightarrow 3,238$ $14 \cdot \sqrt{(1-3,142)^2 + (2-5,428)^2} \Rightarrow 4,043$
 $11 \cdot \sqrt{(1-3,142)^2 + (2-5,428)^2} \Rightarrow 4,043$

ir, $k=3$

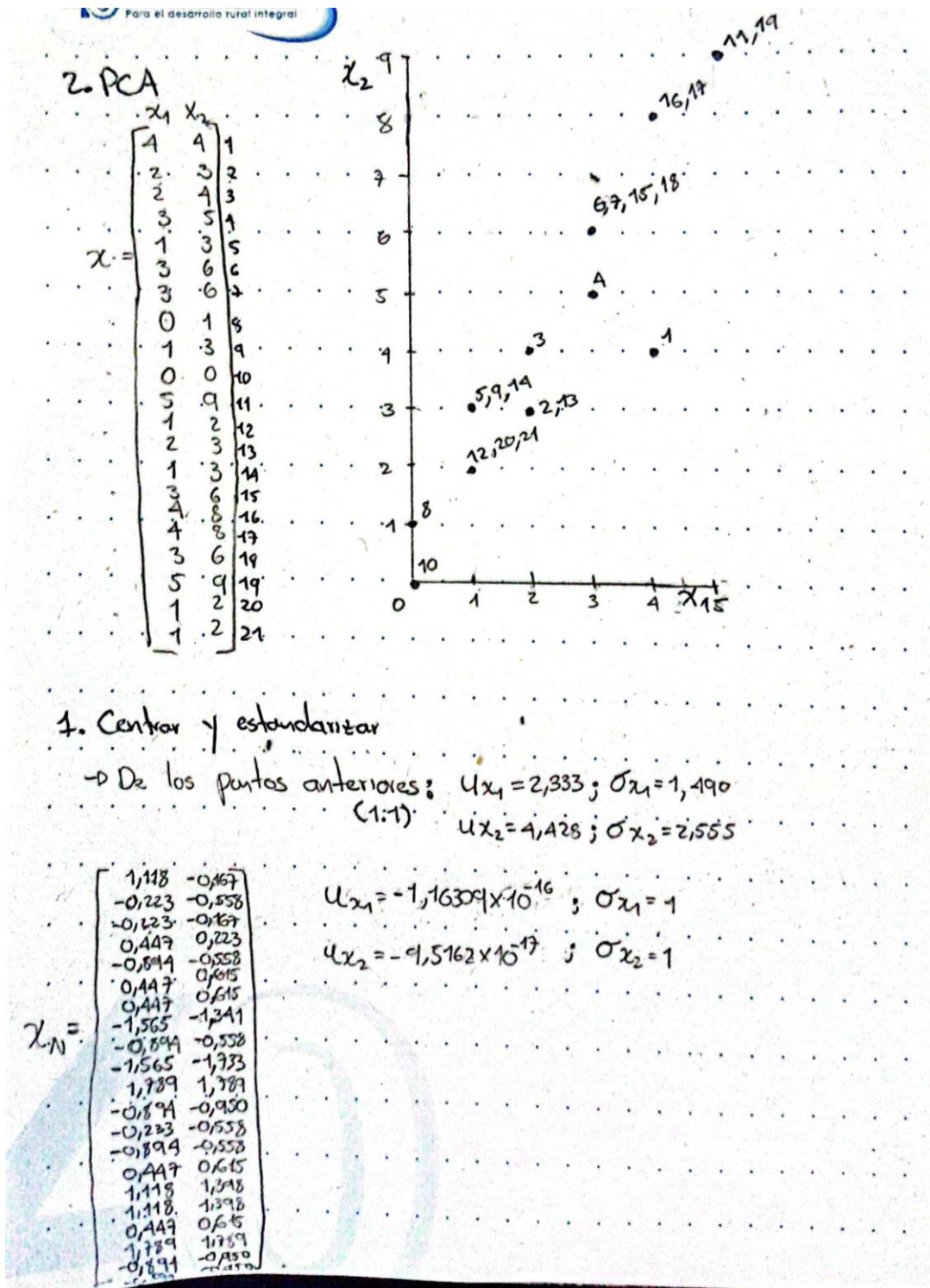
(clústeres).

$dis-C_2 = 1 \sqrt{(1-1,875)^2 + (1-1)^2} \Rightarrow 2,125$	$11 \sqrt{(5-1,875)^2 + (1-1)^2} \Rightarrow 5,896$
$2 \sqrt{(2-1,875)^2 + (3-1)^2} \Rightarrow 1,009$	$12 \sqrt{(1-1,875)^2 + (2-1)^2} \Rightarrow 2,183$
$3 \sqrt{(2-1,875)^2 + (1-1)^2} \Rightarrow 0,125$	$13 \sqrt{(2-1,875)^2 + (3-1)^2} \Rightarrow 1,009$
$4 \sqrt{(3-1,875)^2 + (5-1)^2} \Rightarrow 1,505$	$14 \sqrt{(1-1,875)^2 + (3-1)^2} \Rightarrow 1,328$
$5 \sqrt{(1-1,875)^2 + (3-1)^2} \Rightarrow 1,328$	$15 \sqrt{(3-1,875)^2 + (6-1)^2} \Rightarrow 2,294$
$6 \sqrt{(3-1,875)^2 + (6-1)^2} \Rightarrow 2,294$	$16 \sqrt{(4-1,875)^2 + (8-1)^2} \Rightarrow 4,529$
$7 \sqrt{(3-1,875)^2 + (6-1)^2} \Rightarrow 2,294$	$17 \sqrt{(4-1,875)^2 + (8-1)^2} \Rightarrow 4,529$
$8 \sqrt{(0-1,875)^2 + (1-1)^2} \Rightarrow 3,537$	$18 \sqrt{(3-1,875)^2 + (6-1)^2} \Rightarrow 2,294$
$9 \sqrt{(1-1,875)^2 + (3-1)^2} \Rightarrow 1,328$	$19 \sqrt{(5-1,875)^2 + (9-1)^2} \Rightarrow 5,896$
$10 \sqrt{(0-1,875)^2 + (0-1)^2} \Rightarrow 4,417$	$20 \sqrt{(1-1,875)^2 + (2-1)^2} \Rightarrow 2,183$
	$21 \sqrt{(1-1,875)^2 + (2-1)^2} \Rightarrow 2,183$
$dis-C_1 = 1 \sqrt{(4-2)^2 + (4-3,833)^2} \Rightarrow 2,006$	$11 \sqrt{(5-2)^2 + (9-3,833)^2} \Rightarrow 5,974$
$2 \sqrt{(2-2)^2 + (3-3,833)^2} \Rightarrow 0,833$	$12 \sqrt{(1-2)^2 + (2-3,833)^2} \Rightarrow 2,088$
$3 \sqrt{(2-2)^2 + (4-3,833)^2} \Rightarrow 0,166$	$13 \sqrt{(2-2)^2 + (3-3,833)^2} \Rightarrow 0,833$
$4 \sqrt{(3-2)^2 + (5-3,833)^2} \Rightarrow 1,536$	$14 \sqrt{(1-2)^2 + (3-3,833)^2} \Rightarrow 1,301$
$5 \sqrt{(1-2)^2 + (3-3,833)^2} \Rightarrow 1,301$	$15 \sqrt{(3-2)^2 + (6-3,833)^2} \Rightarrow 2,386$
$6 \sqrt{(3-2)^2 + (6-3,833)^2} \Rightarrow 2,386$	$16 \sqrt{(4-2)^2 + (8-3,833)^2} \Rightarrow 4,621$
$7 \sqrt{(3-2)^2 + (6-3,833)^2} \Rightarrow 2,386$	$17 \sqrt{(4-2)^2 + (8-3,833)^2} \Rightarrow 4,621$
$8 \sqrt{(0-2)^2 + (1-3,833)^2} \Rightarrow 3,468$	$18 \sqrt{(3-2)^2 + (6-3,833)^2} \Rightarrow 2,386$
$9 \sqrt{(1-2)^2 + (3-3,833)^2} \Rightarrow 1,301$	$19 \sqrt{(5-2)^2 + (9-3,833)^2} \Rightarrow 5,974$
$10 \sqrt{(0-2)^2 + (0-3,833)^2} \Rightarrow 4,323$	$20 \sqrt{(1-2)^2 + (2-3,833)^2} \Rightarrow 2,088$
	$21 \sqrt{(1-2)^2 + (2-3,833)^2} \Rightarrow 2,088$

Siguiendo la lógica con cada iteración hasta que no haya cambios en las etiquetas tenemos:

	x1	x2	Inicio	centroide x1	centroide x2	dis_c3_3	dis_c2_2	dis_c1_1	New_labels	centroide x1_it2	centroide x2_it2	dis_c3_3_it2	dis_c2_2_it2	dis_c1_1_it2	New_labels	centroide x1_it3	centroide x2_it3	dis_c3_3_it3	dis_c2_2_it3	dis_c1_1_it3	New_labels	
1		4	4	3	3,14285714	5,428571429	1,66598626	2,125	2,00693243	3	3,7	6,7	2,71661554	2	3,49857114	2	3,333333333	6,444444444	2,533723167	1,414213562	3,816084381	2
2		2	3	3			2,68404203	1,00778222	0,833333333	1			4,07185461	1	1,28062485	2		3,693504475		1,600781059	2	
3		2	4	3			1,82946407	0,125	0,166666667	2			3,19061123	0	2,05912603	2		2,784436464	0,632455532	2,358495283	2	
4		3	5	3			0,45175395	1,50519932	1,53659074	3			1,83847763	1,41421356	3,44093011	2		1,482407118	1,264911064	3,75	2	
5		1	3	2	1,875	4	3,23879544	1,32876823	1,30170828	1	2	4	4,580393	1,41421356	0,8	1	2,6	3,8	4,160365606	1,788854382	1,030776406	1
6		3	6	2			0,58901509	2,29469497	2,38630351	3			0,98994949	2,23606798	4,29418211	3		0,555555556	2,236067977	4,589389938	3	
7		3	6	2			0,58901509	2,29469497	2,38630351	3			0,98994949	2,23606798	4,29418211	3		0,555555556	2,236067977	4,589389938	3	
8		0	1	2			5,4304508	3,53774292	3,46810867	1			6,7955868	3,60555128	1,56204994	1		6,38381441	3,820994635	1,25	1	
9		1	3	2			3,23879544	1,32876823	1,30170828	1			4,580393	1,41421356	0,8	1		4,160365606	1,788854382	1,030776406	1	
10	0	0	0	1	2	3,833333333	6,27271383	4,41764926	4,32370726	1	1	2,2	7,65375725	4,47213595	2,41660919	1	0,75	2	7,255478985	4,604345773	2,136000936	1
11		5	9	1			4,02542937	5,89623821	5,97448278	3			2,64196896	5,83095189	7,88923317	3		3,051006715	5,727128425	8,189169677	3	
12		1	2	1			4,04313477	2,18303115	2,08832735	1			5,42033209	2,23606798	0,2	1		5,019714221	2,408318916	0,25	1	
13		2	3	1			2,68404203	1,00778222	0,833333333	1			4,07185461	1	1,28062485	2		3,693504475		1,600781059	2	
14		1	3	1			3,23879544	1,32876823	1,30170828	1			4,580393	1,41421356	0,8	1		4,160365606	1,788854382	1,030776406	1	
15		3	6	3			0,58901509	2,29469497	2,38630351	3			0,98994949	2,23606798	4,29418211	3		0,555555556	2,236067977	4,589389938	3	
16		4	8	3			2,71052371	4,52941773	4,62180821	3			1,33416641	4,47213595	6,52993109	3		1,692394024	4,427188724	6,823672032	3	
17		4	8	3			2,71052371	4,52941773	4,62180821	3			1,33416641	4,47213595	6,52993109	3		1,692394024	4,427188724	6,823672032	3	
18		3	6	1			0,58901509	2,29469497	2,38630351	3			0,98994949	2,23606798	4,29418211	3		0,555555556	2,236067977	4,589389938	3	
19		5	9	2			4,02542937	5,89623821	5,97448278	3			2,64196896	5,83095189	7,88923317	3		3,051006715	5,727128425	8,189169677	3	
20		1	2	2			4,04313477	2,18303115	2,08832735	1			5,42033209	2,23606798	0,2	1		5,019714221	2,408318916	0,25	1	
21		1	2	2			4,04313477	2,18303115	2,08832735	1			5,42033209	2,23606798	0,2	1		5,019714221	2,408318916	0,25	1	

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.Cuál es la matriz de covarianza

2.2. Cuáles son los eigenvalores

2. Calcular la matriz de covarianza

$$\Sigma = \begin{bmatrix} \sigma_{x_1}^2 & \text{Cov}(x_1, x_2) \\ \text{Cov}(x_2, x_1) & \sigma_{x_2}^2 \end{bmatrix} = \begin{bmatrix} 1 & 0,9375 \\ 0,9375 & 1 \end{bmatrix}$$

3. Hallar eigenvalores

$$\det \left(\begin{bmatrix} 1 & 0,9375 \\ 0,9375 & 1 \end{bmatrix} - \lambda \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \right) = \det \begin{pmatrix} 1-\lambda & 0,9375 \\ 0,9375 & 1-\lambda \end{pmatrix}$$
$$\Rightarrow (1-\lambda)^2 - 0,9375^2 \Rightarrow 1 - 2\lambda + \lambda^2 - 0,878$$
$$\Rightarrow 0,122 - 2\lambda + \lambda^2$$
$$\lambda = \frac{2 \pm \sqrt{4 - 0,488}}{2} \Rightarrow \lambda = \frac{2 \pm \sqrt{3,512}}{2}$$
$$\Rightarrow \lambda = \frac{2 \pm 1,873}{2} \quad \text{por } \lambda_1 = \frac{3,872}{2} \Rightarrow 1,936$$
$$\quad \quad \quad \text{y } \lambda_2 = \frac{0,128}{2} \Rightarrow 0,064$$

→ Se escoge $\lambda_1 = 1,936$ que es el valor mayor

• Cuando se hace la conversión se conserva el $\left(\frac{1,936}{2} \Rightarrow 0,96 \right)$
96% de la información

2.3. Cuál es la varianza explicada por el eigenvalor.

Punto 2.3

Varianza explicada $\lambda_1 = \frac{1,936}{2} \Rightarrow \underline{\underline{0,96}}$

Varianza explicada $\lambda_2 = \frac{0,064}{2} \Rightarrow \underline{\underline{0,032}}$

2.4. Cuál es el valor del eigenvector

Punto 2.4

$$\begin{bmatrix} 1 & 0,9375 \\ 0,9375 & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \lambda \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

① $x_1 + 0,9375 x_2 = 1,936 x_1 \Rightarrow 0,9375 x_2 = 0,936 x_1 \Rightarrow x_1 \approx x_2$

② $0,9375 x_1 + x_2 = 1,936 x_2 \Rightarrow 0,9375 x_1 = 0,936 x_2 \Rightarrow x_2 \approx x_1$

Infinitas soluciones

$V_1 = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$

$\|V_1\| = \sqrt{x_1^2 + x_2^2} = 1 \Rightarrow 1 = \sqrt{2x_1^2} \Rightarrow x_1 = \sqrt{0,5} = 0,707$

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

2.5.Cuál es la matriz proyectada.

Prob 2.5

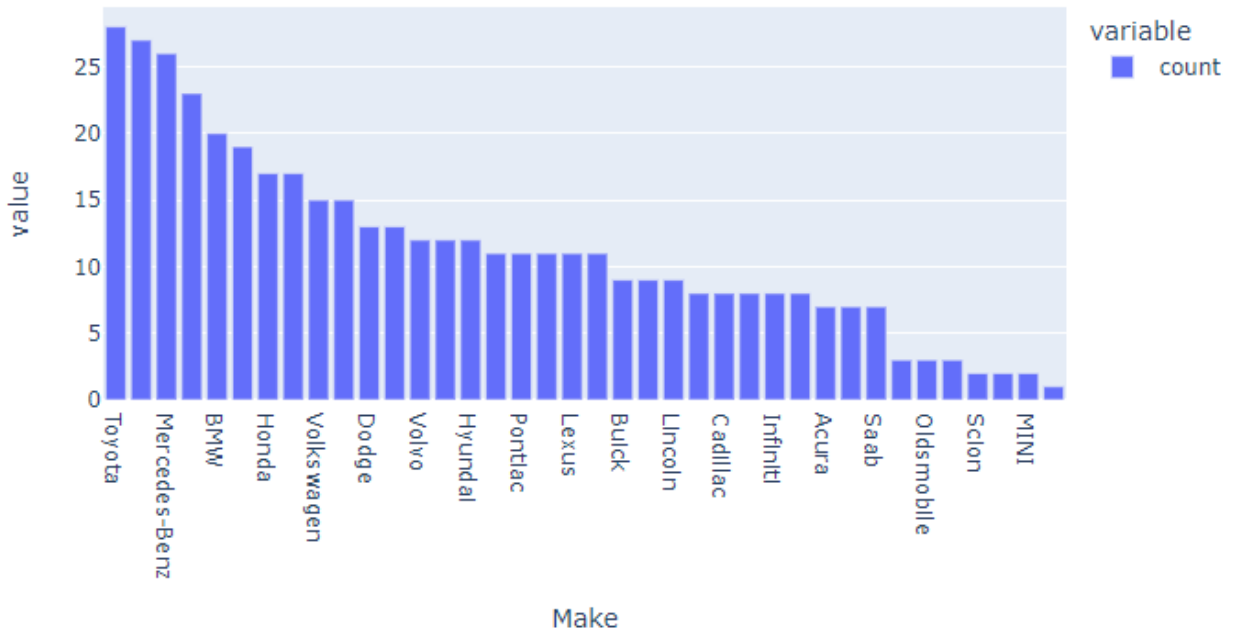
$$\begin{bmatrix}
 1 & 1,118 & -0,167 \\
 2 & -0,223 & -0,558 \\
 3 & -0,223 & -0,167 \\
 4 & 0,447 & 0,223 \\
 5 & -0,894 & -0,558 \\
 6 & 0,447 & 0,615 \\
 7 & 0,447 & 0,615 \\
 8 & -1,565 & -1,341 \\
 9 & -0,894 & -0,558 \\
 10 & -1,565 & -1,789 \\
 11 & 1,789 & 1,789 \\
 12 & -0,894 & -0,950 \\
 13 & -0,223 & -0,558 \\
 14 & -0,894 & -0,558 \\
 & 0,447 & 0,615 \\
 & 1,118 & 1,398 \\
 & 1,118 & 1,398 \\
 & 0,447 & 0,615 \\
 & 1,789 & 1,789 \\
 & -0,894 & -0,950 \\
 & -0,894 & -0,950
 \end{bmatrix}
 \cdot
 \begin{bmatrix}
 0,707 \\
 0,707
 \end{bmatrix}
 =
 \begin{bmatrix}
 0,672 \\
 -0,552 \\
 -0,275 \\
 0,480 \\
 -1,026 \\
 0,750 \\
 0,750 \\
 -2,051 \\
 -1,026 \\
 -2,331 \\
 2,529 \\
 -1,303 \\
 -0,559 \\
 -1,026 \\
 0,750 \\
 1,778 \\
 1,778 \\
 0,750 \\
 2,529 \\
 -1,303 \\
 -1,303
 \end{bmatrix}$$

3. Utilizando el dataset del proyecto data/CARS.csv crear: Utilizar la librería de plotly.

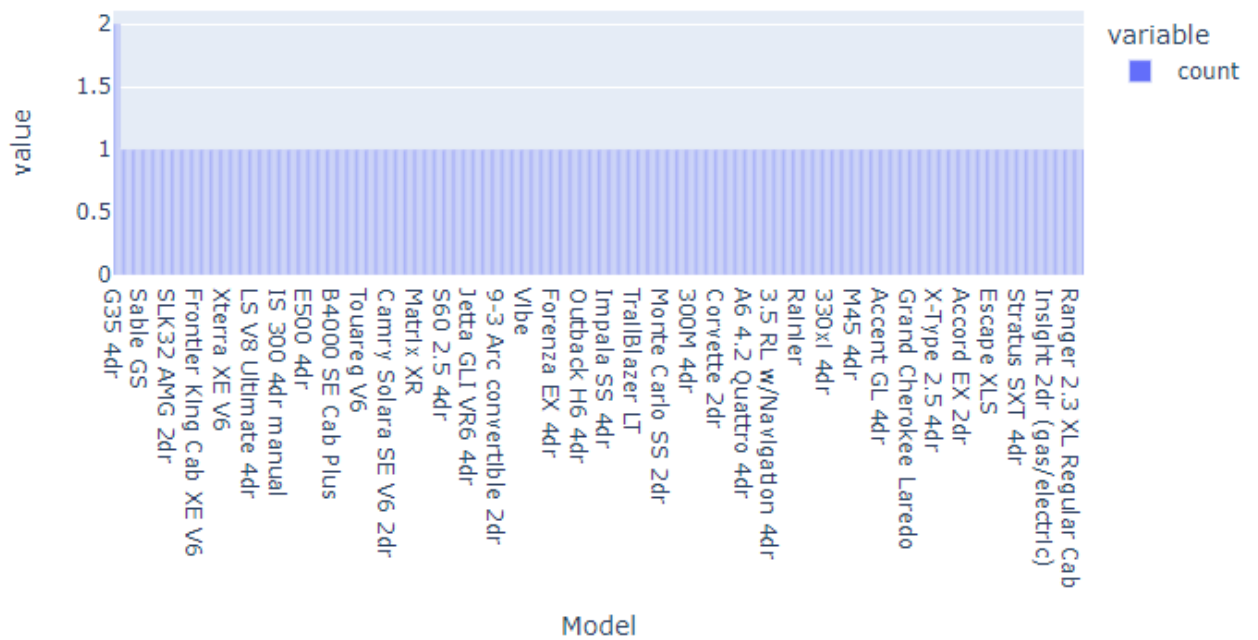
3.1. Distribución de cada variable:

3.1.1. Para las variables categóricas un gráfico de barras. Categoría número de observaciones.

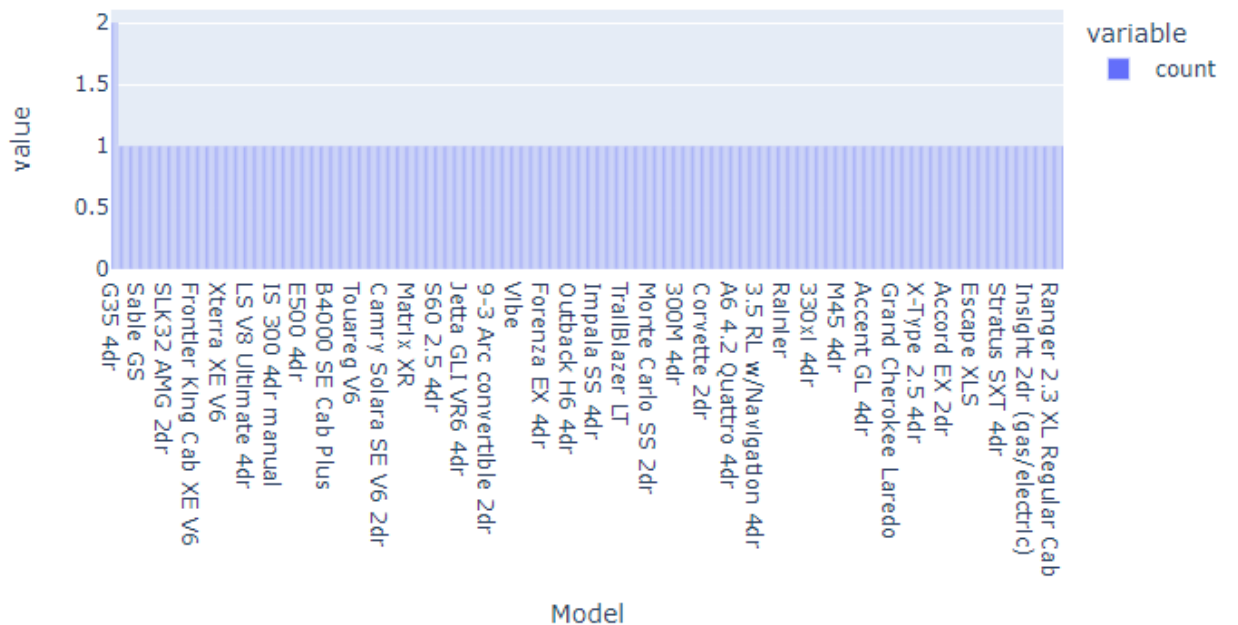
Frecuencia categoria marcas



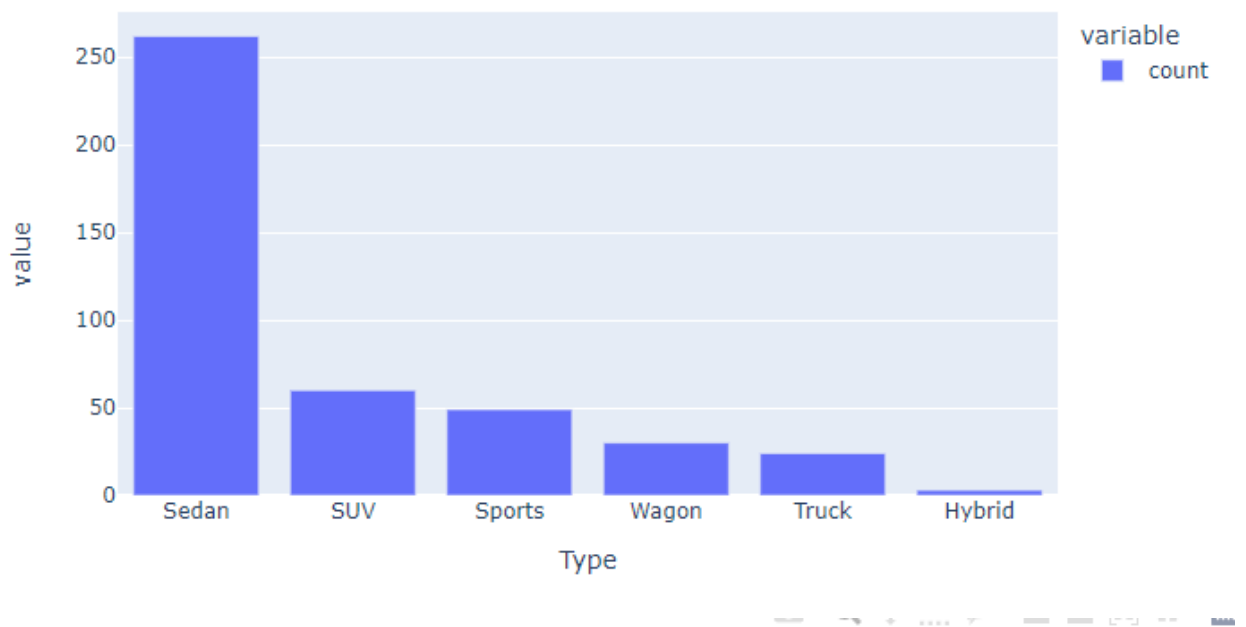
Frecuencia categoria modelo



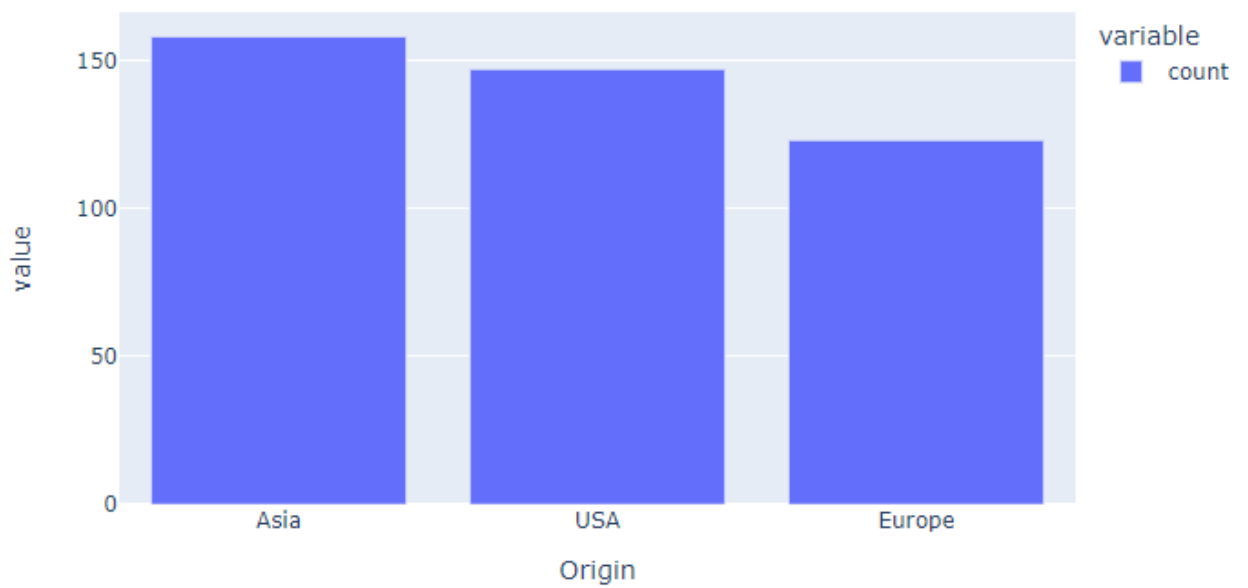
Frecuencia categoria modelo



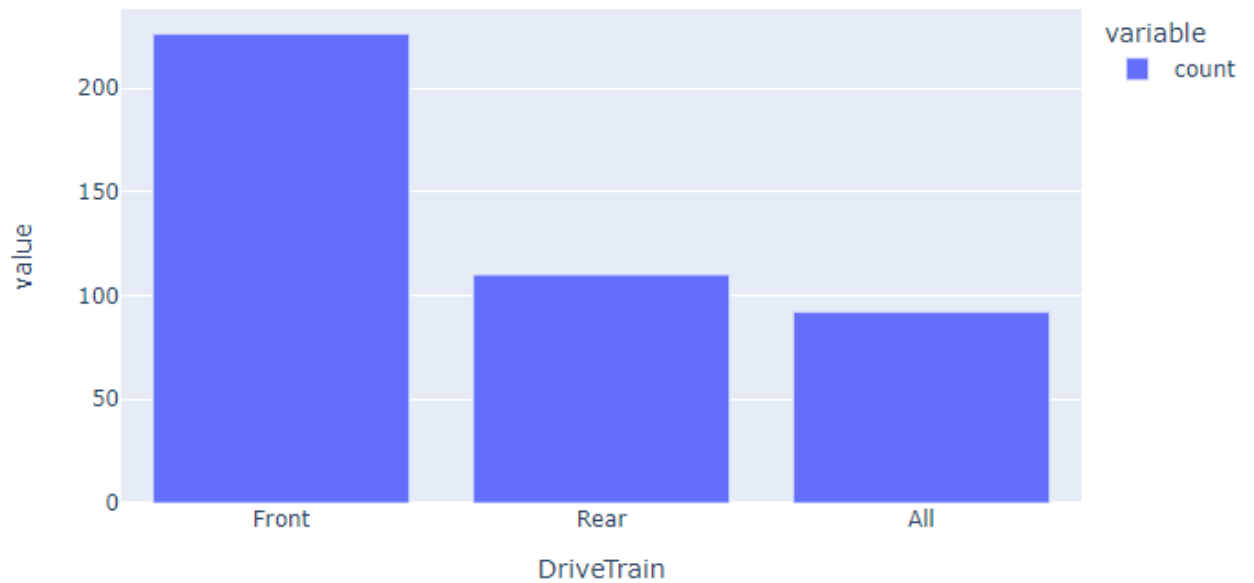
Frecuencia categoria tipo



Frecuencia categoria origen

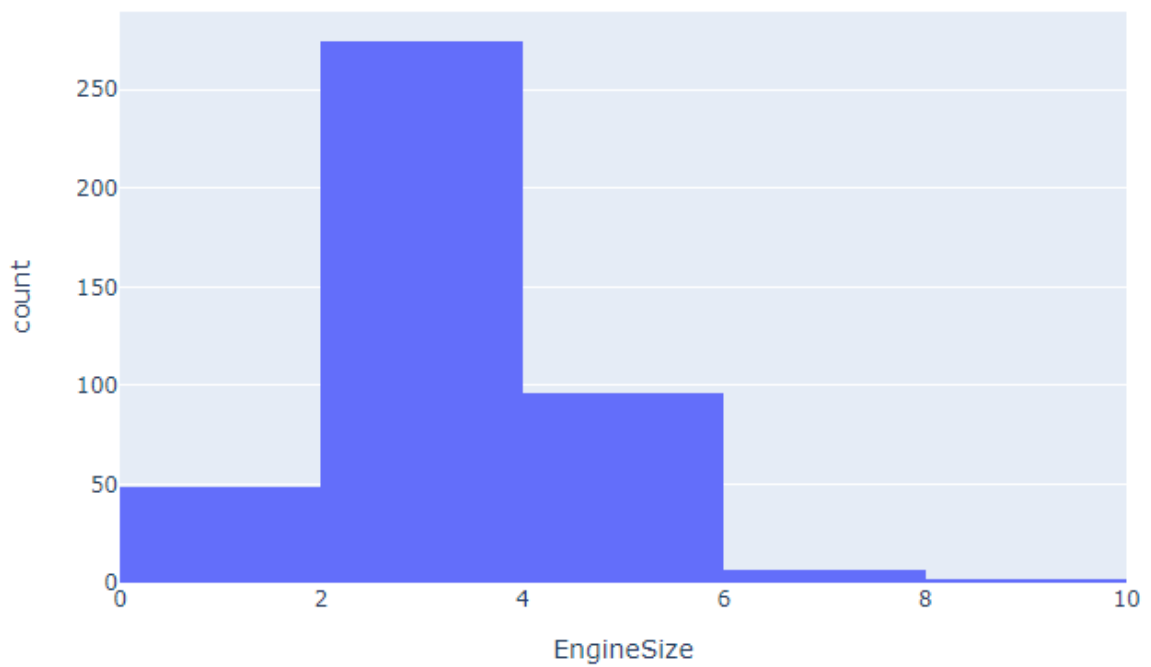
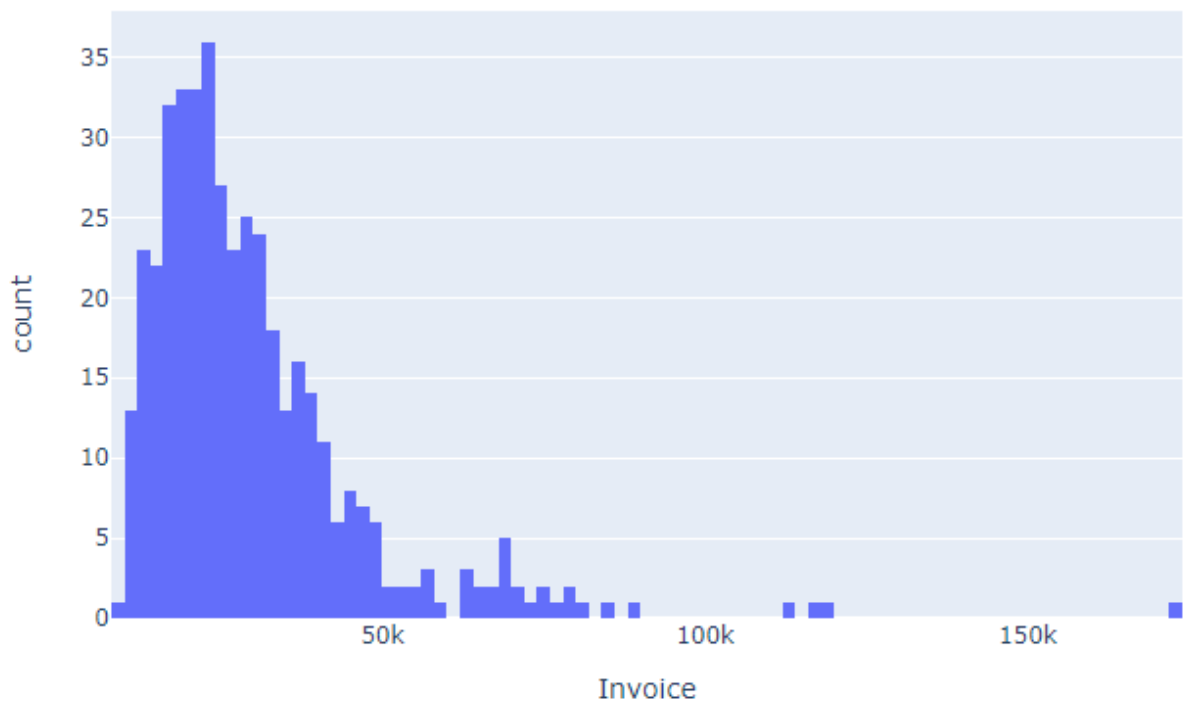


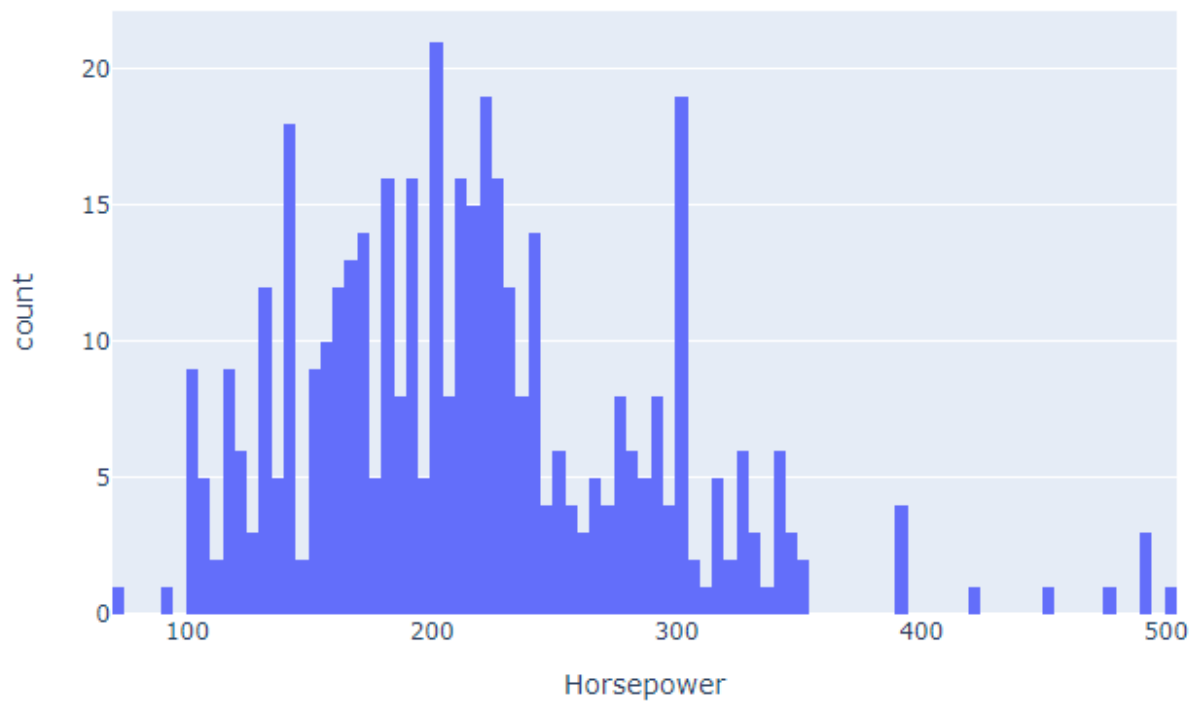
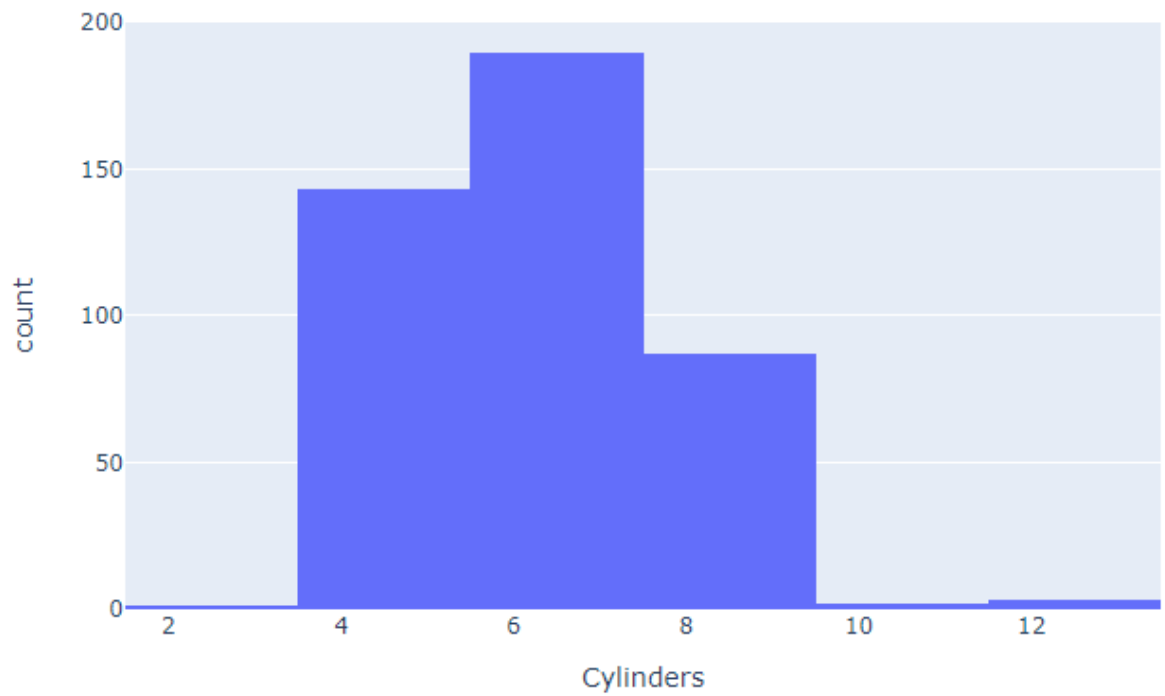
Frecuencia categoria posicion tren traccion

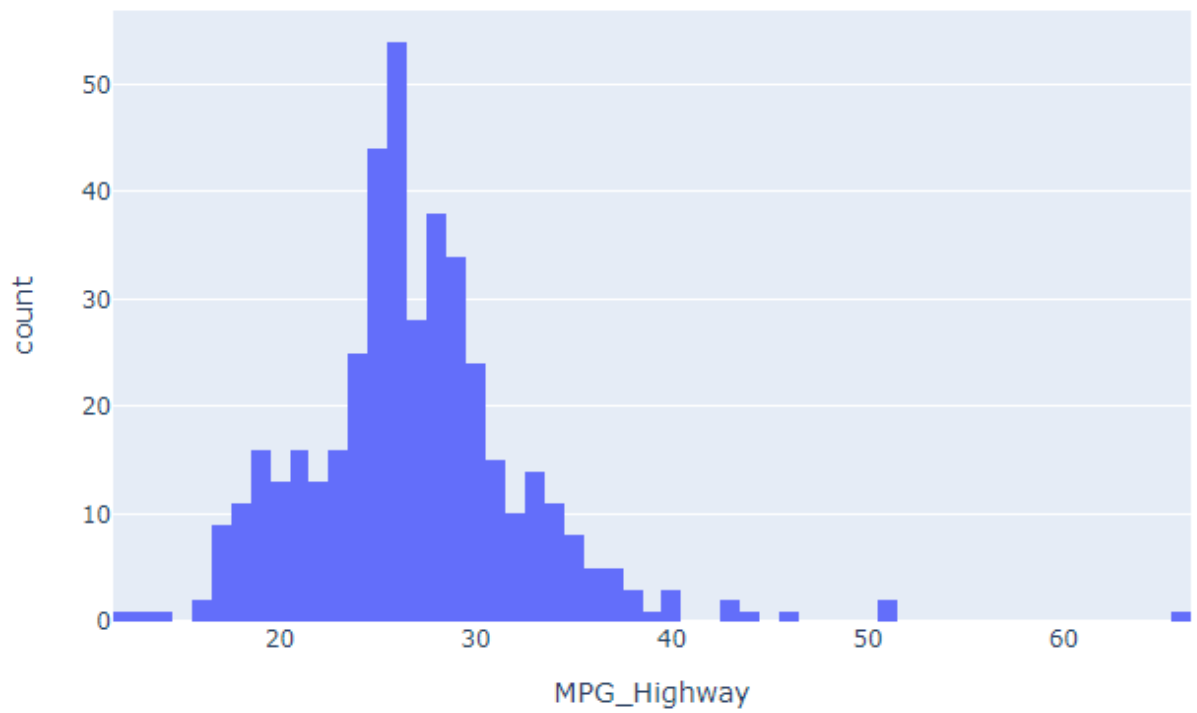
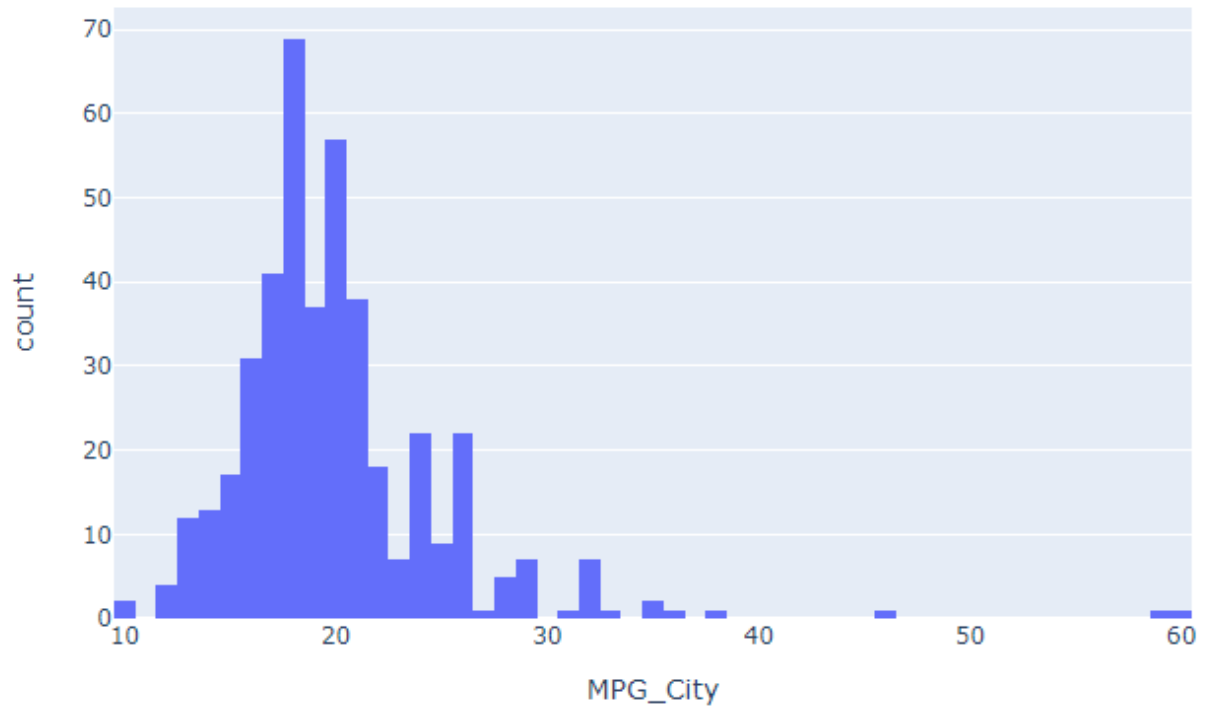


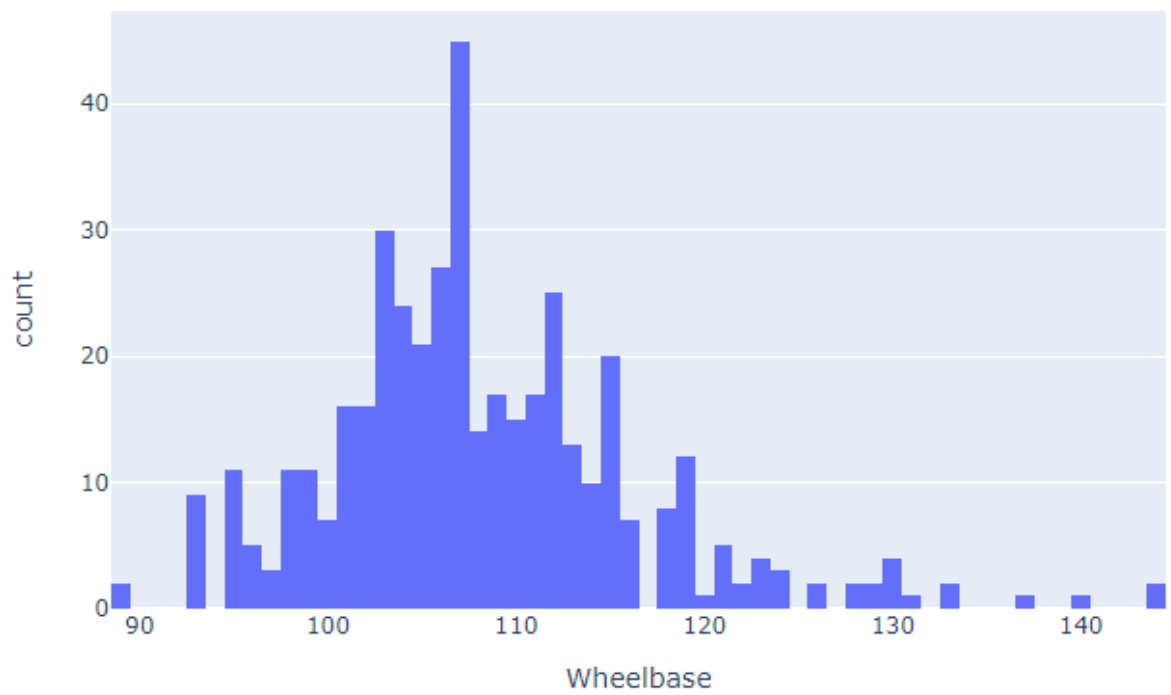
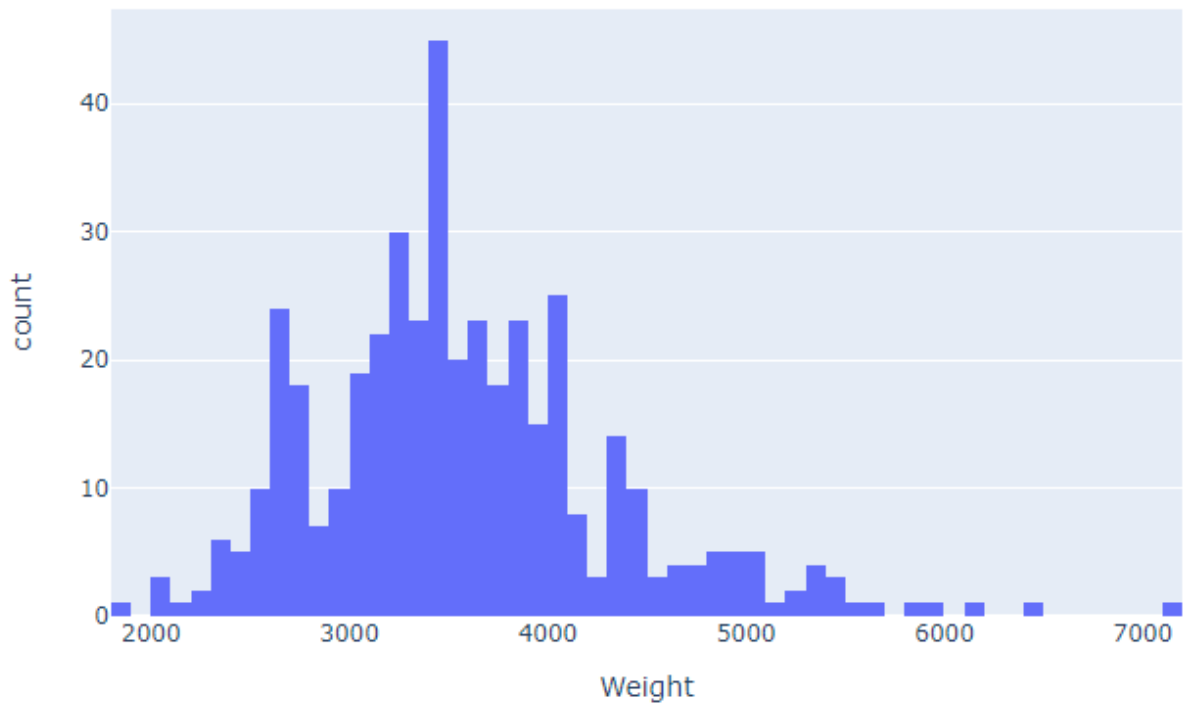
3.1.2. Para las variables numéricas crear histogramas. Listar los modelos de carros que están más lejos de 5 estándares de desviación, y serían considerados outliers. Hacer test de si es una distribución normal o no.

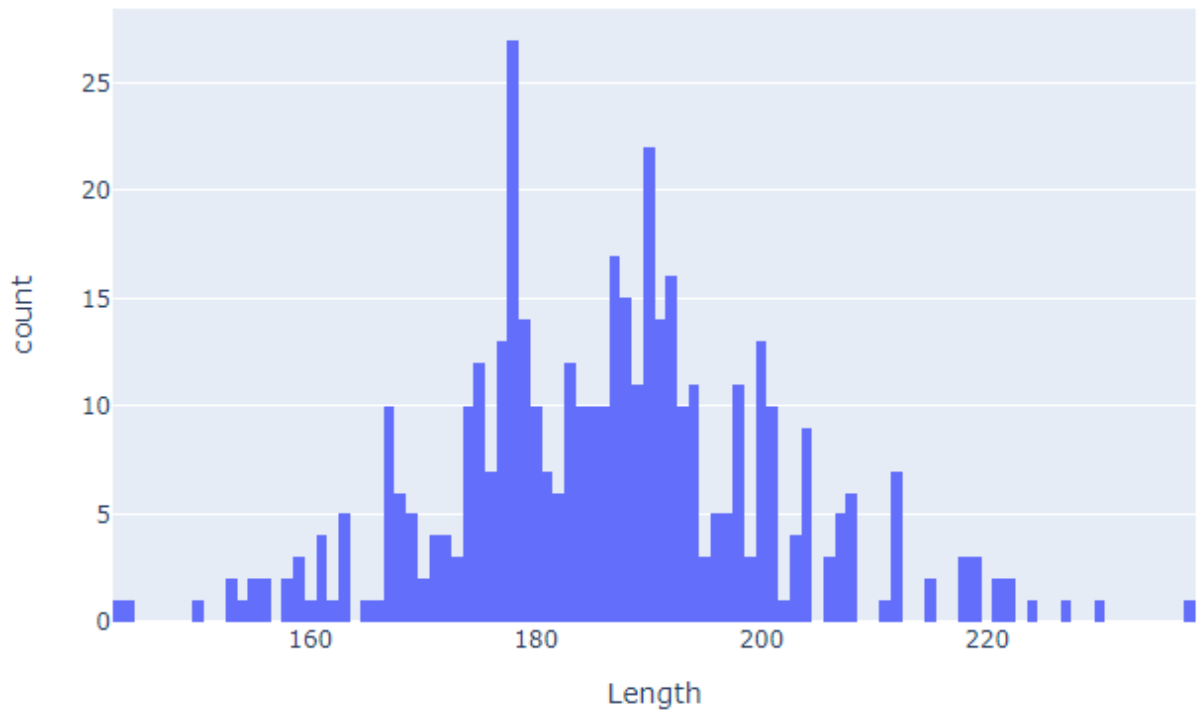
Histogramas:





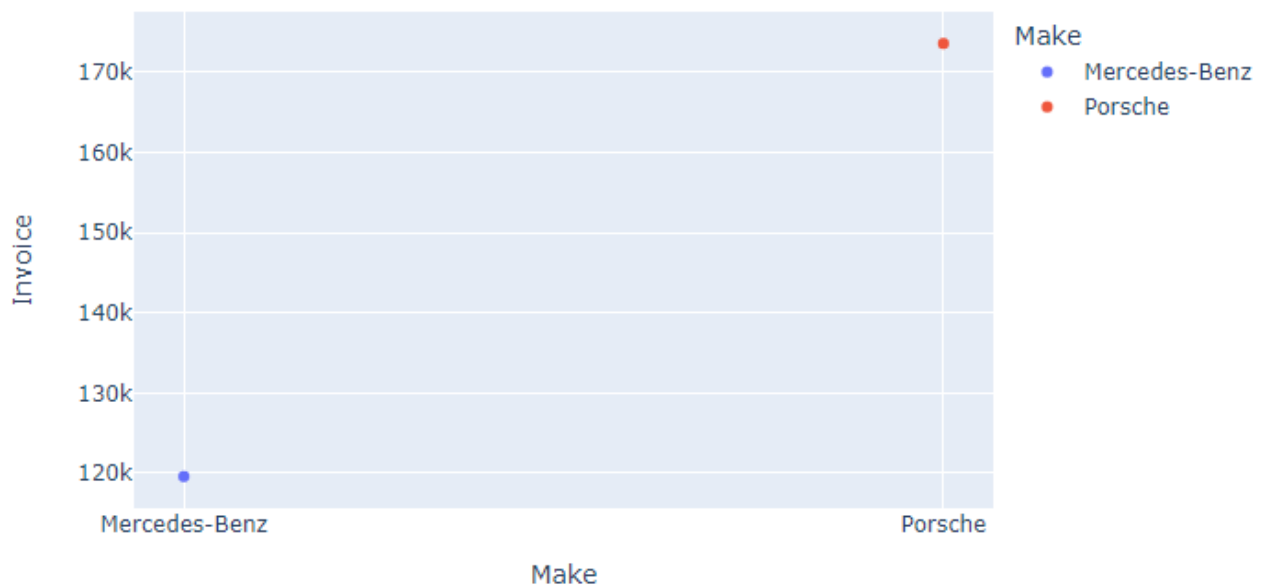




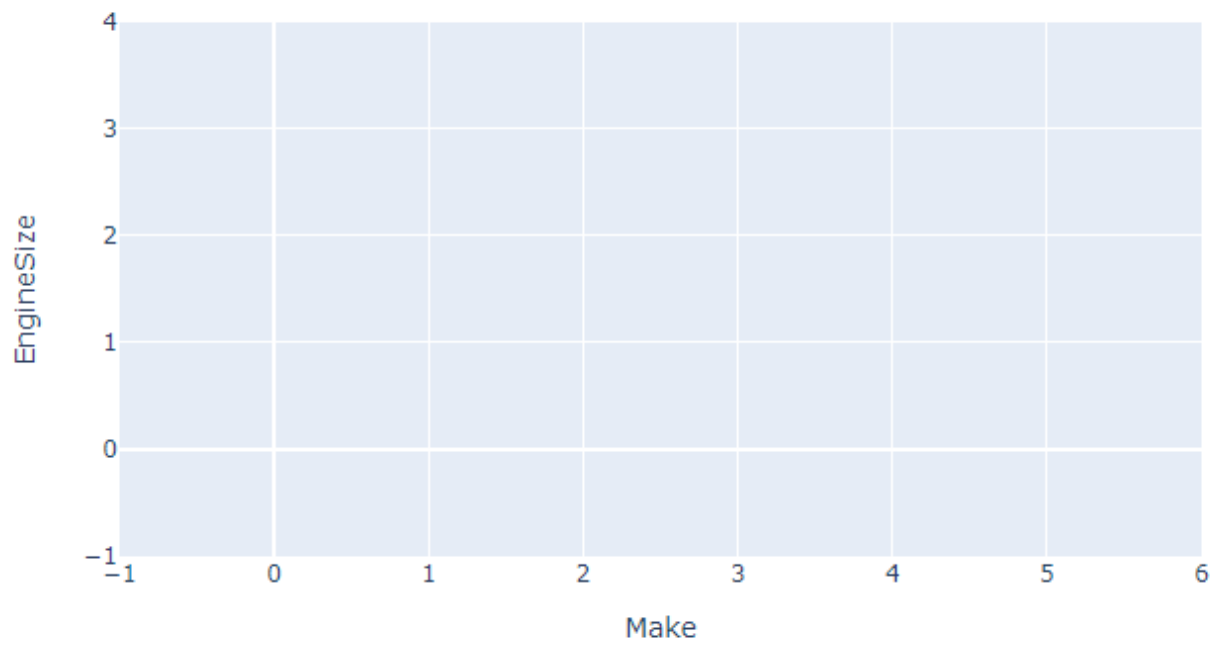


Análisis de Outliers 5 desviaciones estándar de diferencia:

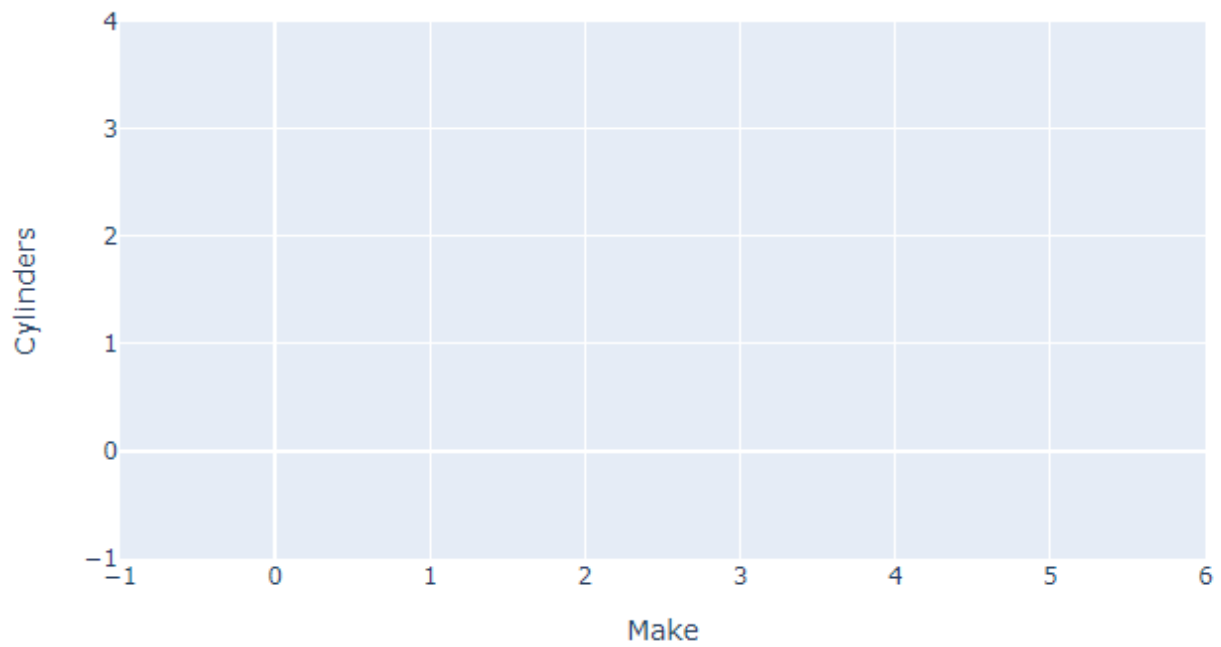
Outliers Invoice



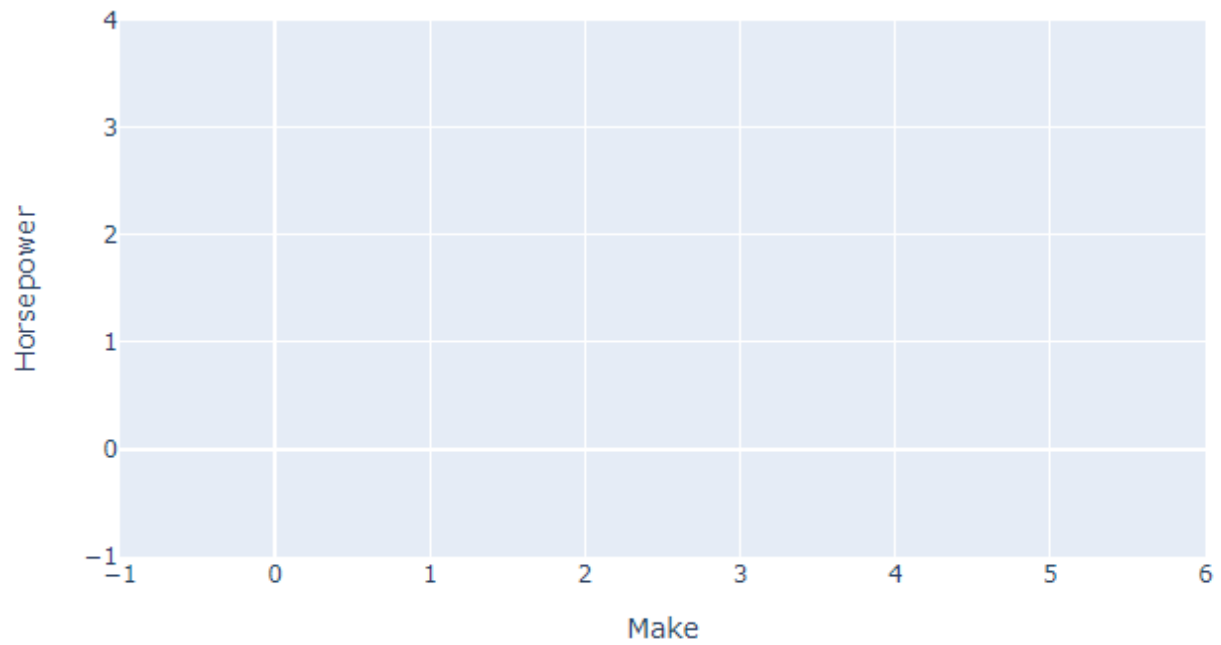
Outliers Engine Size



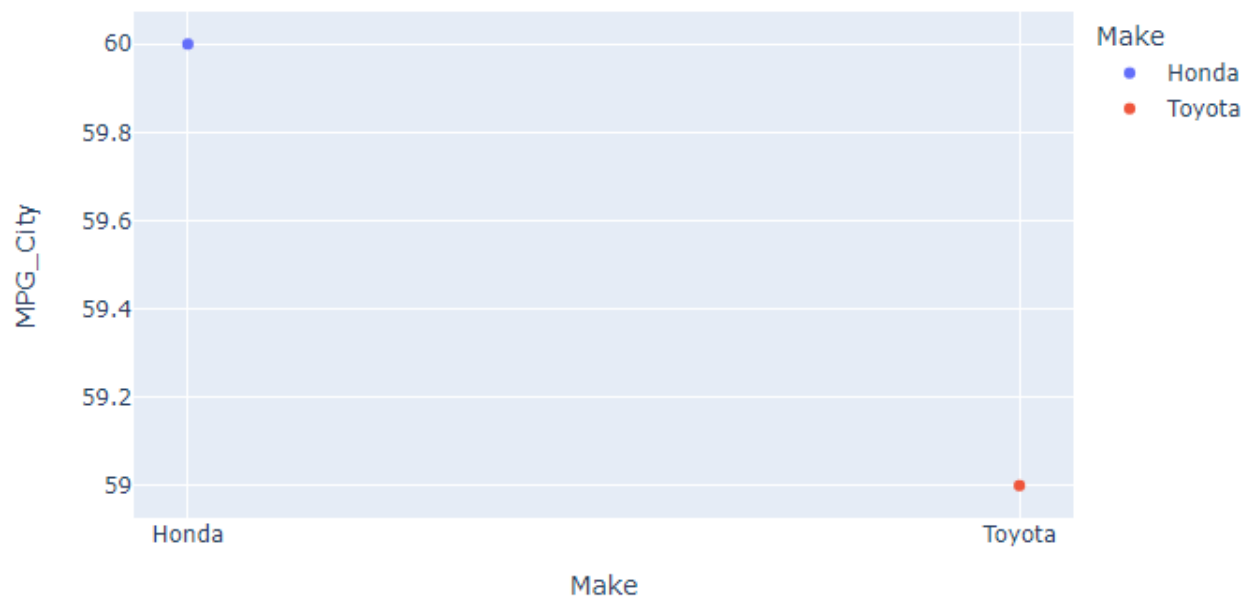
Outliers Cylinders



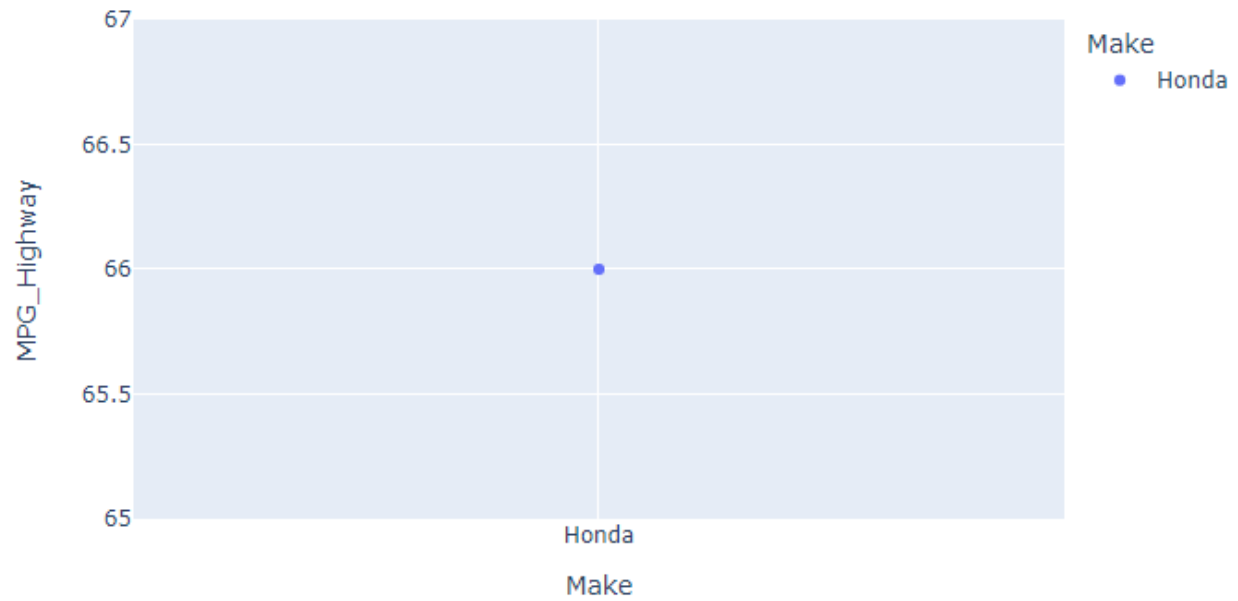
Outliers Horse Power



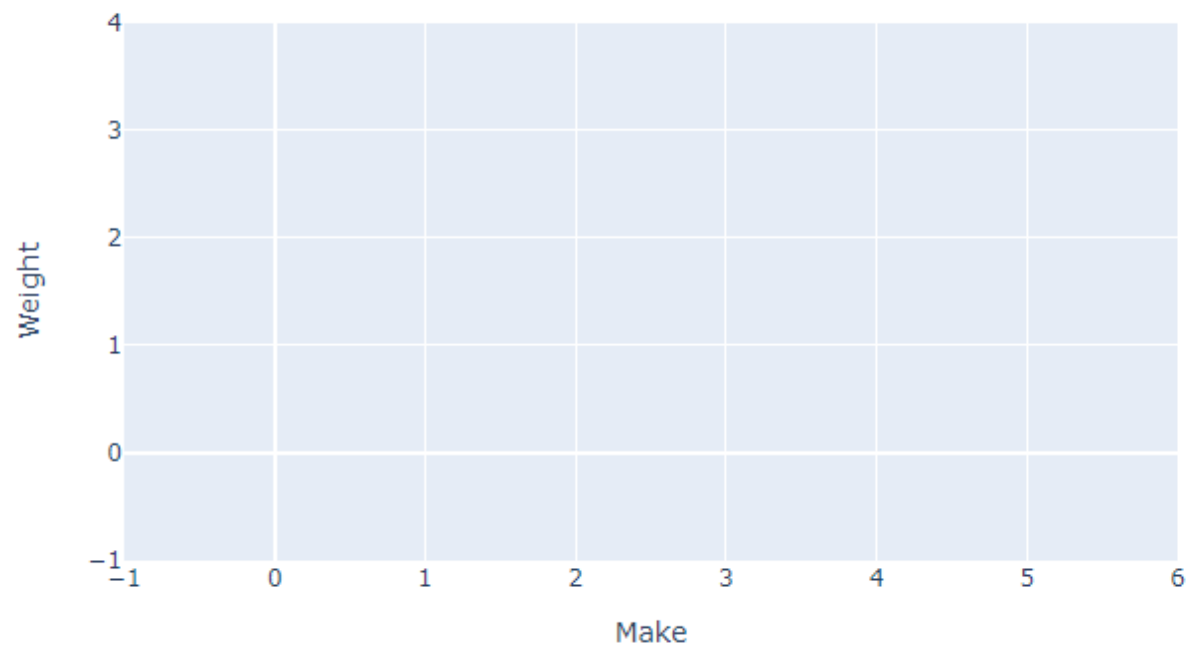
Outliers Miles per Gallon City



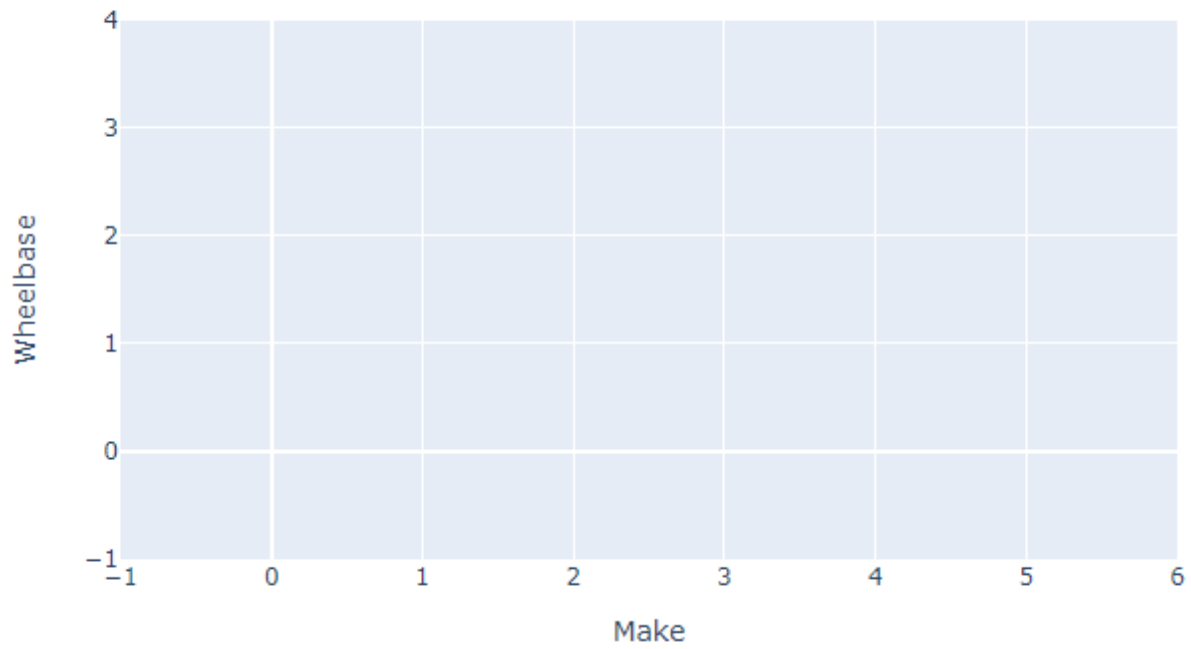
Outliers Miles per Gallon Highway



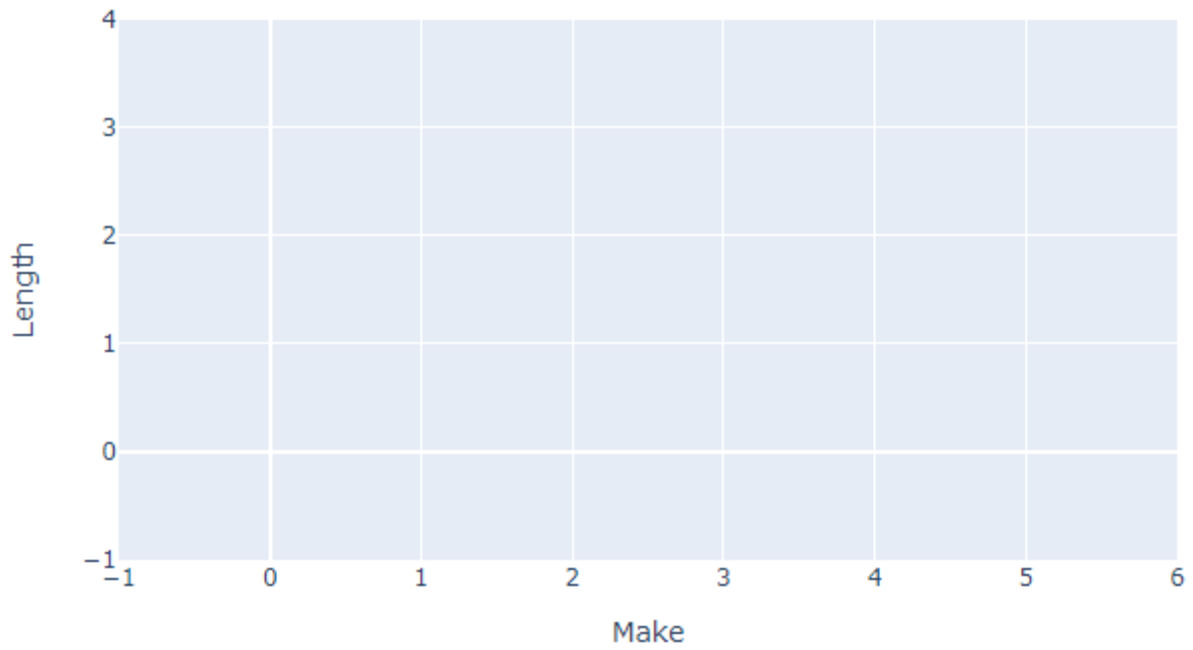
Outliers Weigth



Outliers Wheelbase



Outliers Length



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, cuáles 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 todos los modelos de carro que fueron catalogados como un outlier. (Puede utilizar. `query('Model in["MDX","TSX 4dr"]')`). Existe alguna variación en la correlación.