

▼ Actividad 3 - Transformaciones

Frida Cano Falcón - A01752953

▼ Extracción de datos

```
1 from google.colab import drive
2 drive.mount('/content/drive')
```

Drive already mounted at /content/drive; to attempt to forcibly remount, call drive.mount

```
1 %cd "/content/drive/MyDrive/7mo Semestre/Estadistica"
2 !ls
```

```
/content/drive/MyDrive/Semestres/7mo Semestre/Estadistica
Act1_Distribuciones_FridaCano_A01752953.ipynb
Act2_ExplorandoBases_FridaCano_A01752953.ipynb
Act3_Transformaciones_FridaCanoFalcon_A01752953.ipynb
mc-donalds-menu-1.csv
```

```
1 import pandas as pd
2 import matplotlib.pyplot as plt
3 import numpy as np
4 from scipy import stats
5 # Importing library
6 from scipy.stats import skew, kurtosis
```

```
1 data_base = pd.read_csv('mc-donalds-menu-1.csv')
2 data_base.head()
```

	Category	Item	Serving Size	Calories	Calories from Fat	Total Fat	Total Fat (% Daily Value)	Saturated Fat	Saturated Fat (% Daily Value)	Tr
0	Breakfast	Egg McMuffin	4.8 oz (136 g)	300	120	13.0	20	5.0	25	
1	Breakfast	Egg White	4.8 oz (135 g)	250	70	8.0	12	3.0	15	

▼ Normalización

La variable a transformar para normalizar es: **Sodium**

3	Breakfast	McMuffin	4.8 oz (136 g)	450	250	28.0	43	10.0	52
---	-----------	----------	----------------	-----	-----	------	----	------	----

```
1 data = data_base["Sodium"]
2 # Debido a que el dataframe contiene valores igual a 0 les sumamos 1 a todos los valores.
3 data = data+1
4 print(data)
```

```
0    751
1    771
2    781
3    861
4    881
```

```
...
255   281
256   381
257   191
258   401
259   201
```

Name: Sodium, Length: 260, dtype: int64

```
1 # Verificar si hay valores negativos en el DataFrame
2 hay_ceros = (data <= 0).any().any()
3 if hay_ceros:
4     print("El DataFrame contiene valores negativos.")
5 else:
6     print("El DataFrame no contiene valores negativos.")
```

El DataFrame no contiene valores negativos.

Transformación con Box-Cox

```
1 # Funcion Box-cox
2 transformed_data, lambda_value = stats.boxcox(data)
3
4 # Imprimir el DataFrame con los datos transformados y el valor de lambda
```

5

Valor de lambda: 0.2103467311625166

Las ecuaciones son:

```
1 x = np.linspace(np.min(data),np.max(data),260)
2 # Ecuación práctica
3 eq_1 = np.sqrt(x+1)
4 # Ecuación precisa
5 eq_2 = (np.power((x+1),lambda_value)-1)/lambda_value
```

Analiza la normalidad de las transformaciones obtenidas con los datos originales. Utiliza como argumento de normalidad:

1. Compara las medidas: Mínimo, máximo, media, mediana, cuartil 1 y cuartil 3, sesgo y curtosis.

```
1 datos = ['Datos originales', 'Ec. aproximada', 'Ec. exacta']
2 # Mínimos
3 eq1_min = np.min(eq_1)
4 eq2_min = np.min(eq_2)
5 data_min = np.min(data)
6 mins = [data_min, eq1_min, eq2_min]
7
8 # Máximos
9 eq1_max = np.max(eq_1)
10 eq2_max = np.max(eq_2)
11 data_max = np.max(data)
12 maxs = [data_max, eq1_max, eq2_max]
13
14 # Media
15 eq1_media = np.mean(eq_1)
16 eq2_media = np.mean(eq_2)
17 data_media = np.mean(data)
18 means = [data_media, eq1_media, eq2_media]
19
20 # Mediana
21 eq1_median = np.median(eq_1)
22 eq2_median = np.median(eq_2)
23 data_median = np.median(data)
24 medians = [data_median, eq1_median, eq2_median]
25
26 # Cuartiles
27 eq1_cuartiles = np.percentile(eq_1, [25, 50, 75])
28 eq2_cuartiles = np.percentile(eq_2, [25, 50, 75])
29 data_cuartiles = np.percentile(data, [25, 50, 75])
30 q1s = [data_cuartiles[0], eq1_cuartiles[0], eq2_cuartiles[0]]
31 q3s = [data_cuartiles[2], eq1_cuartiles[2], eq2_cuartiles[2]]
32
```

```

33 # Sesgo
34 eq1_sesgo = stats.skew(eq_1)
35 eq2_sesgo = stats.skew(eq_2)
36 data_sesgo = stats.skew(data)
37 sesgos = [data_sesgo, eq1_sesgo, eq2_sesgo]
38
39 # Curtosis
40 eq1_curt = stats.kurtosis(eq_1)
41 eq2_curt = stats.kurtosis(eq_2)
42 data_curt = stats.kurtosis(data)
43 curts = [data_curt, eq1_curt, eq2_curt]
44
45 # Tabla
46 tabla_info = {'Mínimo': mins, 'Q1': q1s, 'Mediana': medians, 'Media': means, 'Q3': q3s, 'Máximo': maxs}
47 tabla = pd.DataFrame(tabla_info)
48 print(tabla)

```

	Mínimo	Q1	Mediana	Media	Q3	Máximo	\
0	1.000000	108.500000	191.000000	496.750000	866.000000	3601.000000	
1	1.414214	30.033147	42.449892	39.995676	51.980733	60.016664	
2	0.746226	15.138046	18.255121	17.222675	20.301720	21.863709	

	Curtosis
0	2.796412
1	-0.575289
2	1.355971

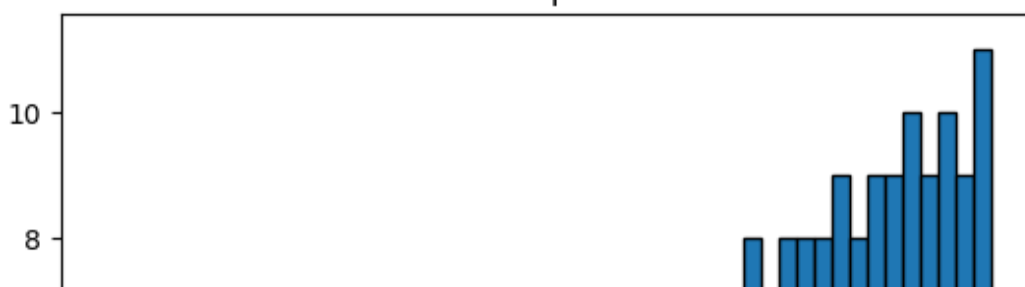
2. Obten el histograma de los 2 modelos obtenidos (exacto y aproximado) y los datos originales.

```

1 # Visualizar el histograma
2 plt.hist(eq_1, 50, edgecolor='black')
3 plt.title("Ecuación aproximada")
4 plt.xlabel("Valores")
5 plt.ylabel("Frecuencia")
6 plt.show()

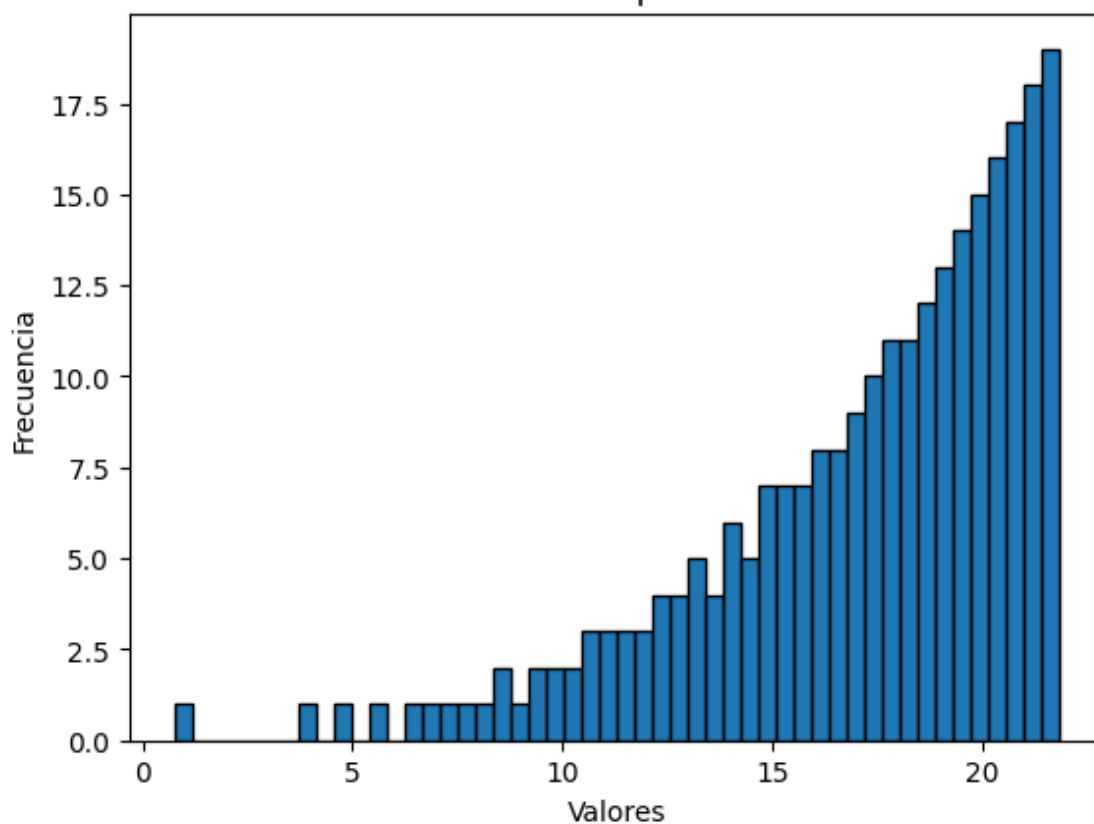
```

Ecuación aproximada

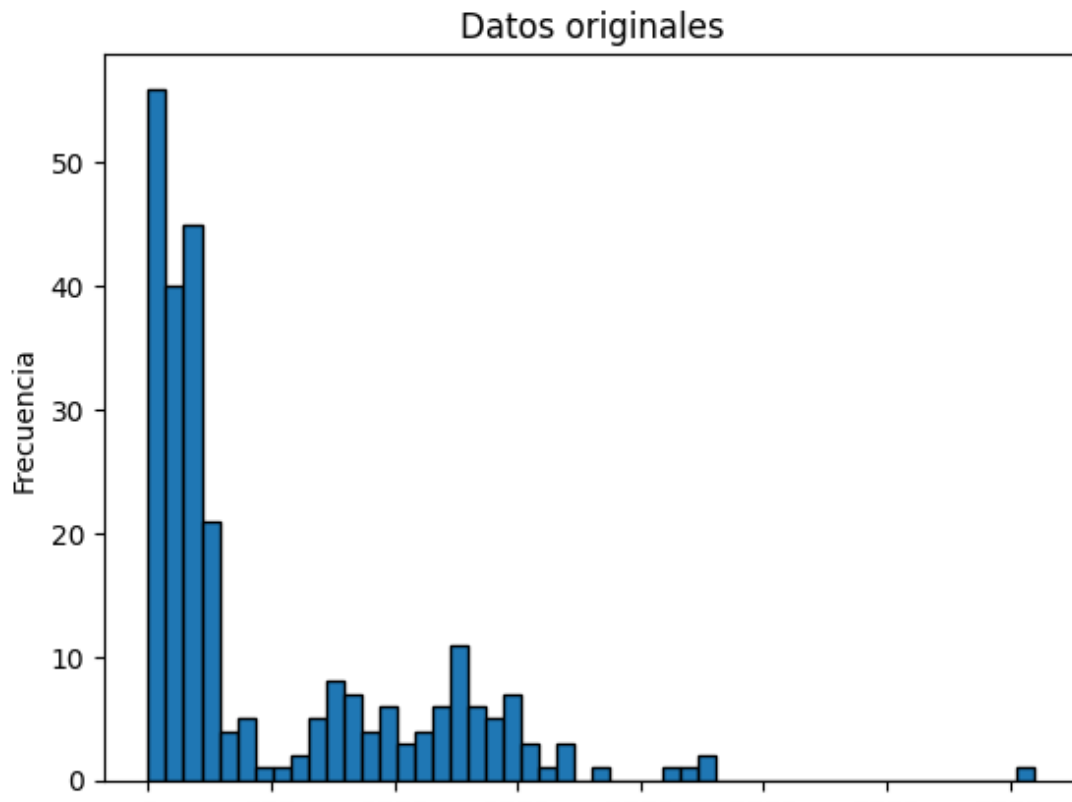


```
1 # Visualizar el histograma
2 plt.hist(eq_2, 50, edgecolor='black')
3 plt.title("Ecuación precisa")
4 plt.xlabel("Valores")
5 plt.ylabel("Frecuencia")
6 plt.show()
```

Ecuación precisa



```
1 # Visualizar el histograma
2 plt.hist(data,50, edgecolor='black')
3 plt.title("Datos originales")
4 plt.xlabel("Valores")
5 plt.ylabel("Frecuencia")
6 plt.show()
```



3. Realiza la prueba de normalidad de Anderson-Darling o de Jarque Bera para los datos transformados y los originales

```

1 _eq1, critical_values_eq1, significance_levels_eq1 = stats.anderson(eq_1, dist='norm')
2 _eq2, critical_values_eq2, significance_levels_eq2 = stats.anderson(eq_2, dist='norm')
3 _data, critical_values_data, significance_levels_data = stats.anderson(data, dist='norm')
4
5 print(f"Resultados para Ecuación Aproximada:")
6 print(f"Estadístico Anderson-Darling: {_eq1}")
7 print("Niveles críticos:", critical_values_eq1)
8 print("Niveles de significancia:", significance_levels_eq1)
9 print("-" * 30)
10
11 print(f"\nResultados para Ecuación Exacto:")
12 print(f"Estadístico Anderson-Darling: {_eq2}")
13 print("Niveles críticos:", critical_values_eq2)
14 print("Niveles de significancia:", significance_levels_eq2)
15 print("-" * 30)
16
17 print(f"\nResultados para Datos originales:")
18 print(f"Estadístico Anderson-Darling: {_data}")
19 print("Niveles críticos:", critical_values_data)
20 print("Niveles de significancia:", significance_levels_data)
21 print("-" * 30)

```

```

Resultados para Ecuación Aproximada:
Estadístico Anderson-Darling: 3.731852911849842
Niveles críticos: [0.567 0.646 0.775 0.904 1.076]

```

Niveles de significancia: [15. 10. 5. 2.5 1.]

Resultados para Ecuación Exacto:

Estadístico Anderson-Darling: 7.267329591457155

Niveles críticos: [0.567 0.646 0.775 0.904 1.076]

Niveles de significancia: [15. 10. 5. 2.5 1.]

Resultados para Datos originales:

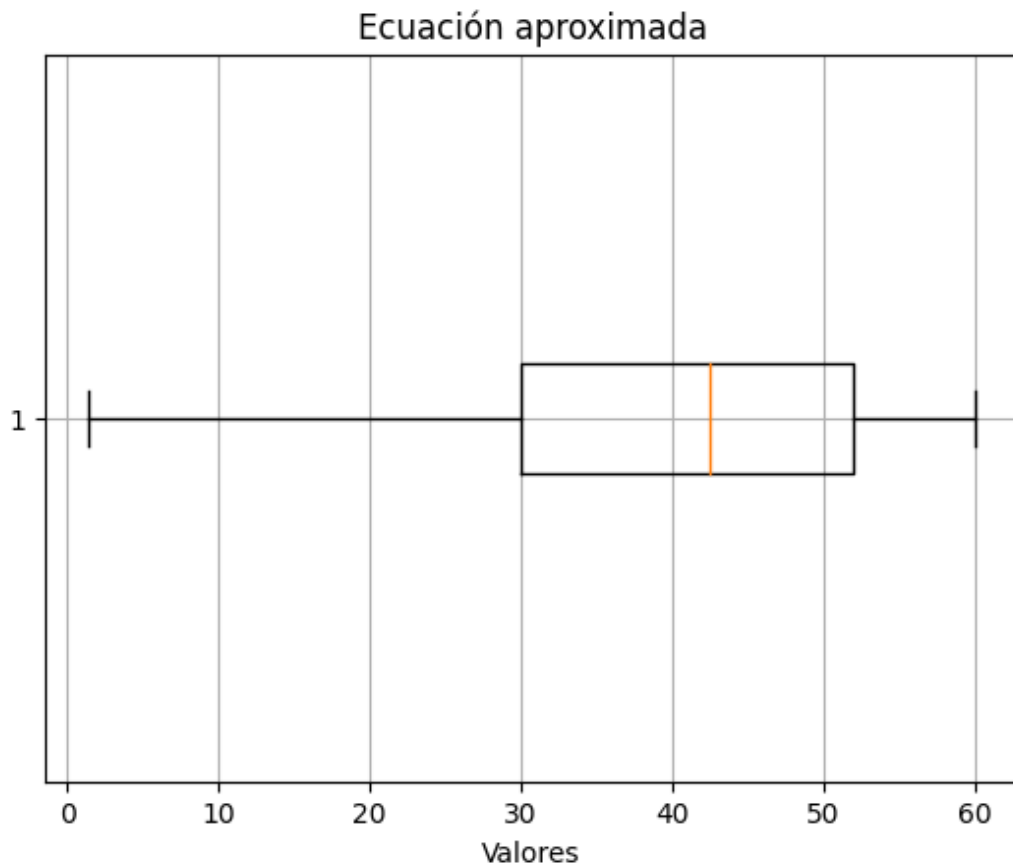
Estadístico Anderson-Darling: 21.405561557715373

Niveles críticos: [0.567 0.646 0.775 0.904 1.076]

Niveles de significancia: [15. 10. 5. 2.5 1.]

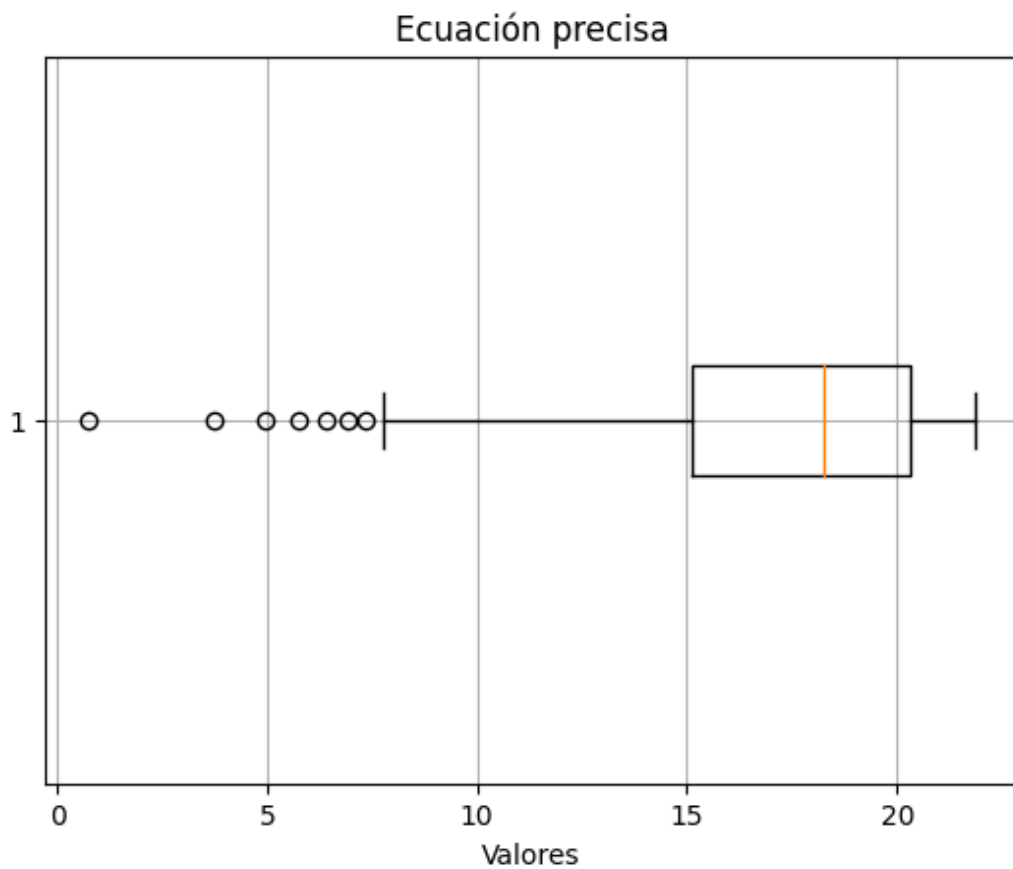
Aplicamos la grafica de bogote para revisar si existen valores atípicos

```
1 # Valores Ecuación aproximada
2 plt.boxplot(eq_1, vert=False)
3 plt.title("Ecuación aproximada")
4 plt.xlabel("Valores")
5 plt.grid(True)
6 plt.show()
```



```
1 # Valores originales
2 plt.boxplot(eq_2,vert=False)
```

```
3 plt.title("Ecuación precisa")
4 plt.xlabel("Valores")
5 plt.grid(True)
6 plt.show()
```



```
1 # Valores originales
2 plt.boxplot(data,vert=False)
3 plt.title("Datos originales")
4 plt.xlabel("Valores")
5 plt.grid(True)
6 plt.show()
```


Datos originales

Detección de anomalías y corrección de base de datos.

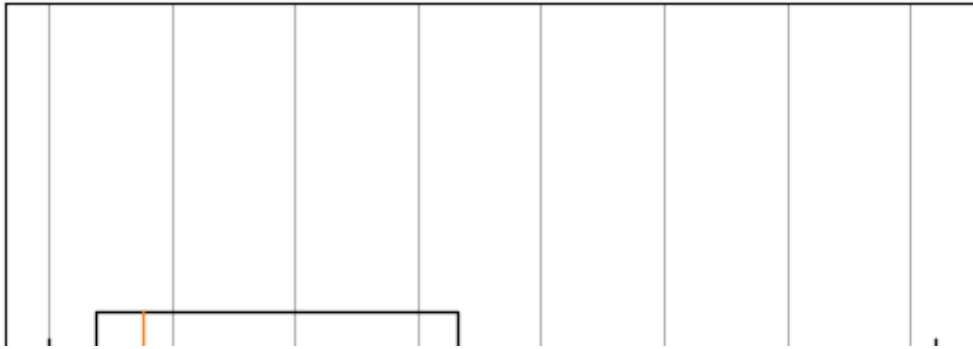
Vemos que despues de aproximadamente 2000 existen valores atípicos, por lo que estos serán eliminados.

```
1 data = data.drop(data[data >=2000].index)
2 print(data)
```

```
0      751
1      771
2      781
3      861
4      881
...
255    281
256    381
257    191
258    401
259    201
Name: Sodium, Length: 255, dtype: int64
```

```
1 # Valores originales
2 plt.boxplot(data,vert=False)
3 plt.title("Datos originales")
4 plt.xlabel("Valores")
5 plt.grid(True)
6 plt.show()
```

Datos originales



Utilizamos la transformación de Yeo Johnson

```
1 # Aplicar la transformación Yeo-Johnson
2 transformed_data, lambda_value_yj = stats.yeojohnson(data)
3
4 print("Valor lambda:", lambda_value_yj)
```

Valor lambda: 0.2041285582397882



Ecuación:

```
1 x = np.linspace(np.min(data), np.max(data), 260)
2
3 # Ecuacion aproximada
4 yj_eq1 = x ** 2
5 # Ecuacion precisa
6 yj_eq2 = (((x+1)**lambda_value_yj) - 1) / lambda_value_yj
```

Análisis:

```
1 datos = ['Datos originales', 'Ec. aproximada', 'Ec. exacta']
2 # Mínimos
3 eq1_min = np.min(yj_eq1)
4 eq2_min = np.min(yj_eq2)
5 data_min = np.min(data)
6 mins = [data_min, eq1_min, eq2_min]
7
8 # Máximos
9 eq1_max = np.max(yj_eq1)
10 eq2_max = np.max(yj_eq2)
11 data_max = np.max(data)
12 maxs = [data_max, eq1_max, eq2_max]
13
14 # Media
15 eq1_media = np.mean(yj_eq1)
16 eq2_media = np.mean(yj_eq2)
```

```

17 data_media = np.mean(data)
18 means = [data_media, eq1_media, eq2_media]
19
20 # Mediana
21 eq1_median = np.median(yj_eq1)
22 eq2_median = np.median(yj_eq2)
23 data_median = np.median(data)
24 medians = [data_median, eq1_median, eq2_median]
25
26 # Cuartiles
27 eq1_cuartiles = np.percentile(yj_eq1, [25, 50, 75])
28 eq2_cuartiles = np.percentile(yj_eq2, [25, 50, 75])
29 data_cuartiles = np.percentile(data, [25, 50, 75])
30 q1s = [data_cuartiles[0], eq1_cuartiles[0], eq2_cuartiles[0]]
31 q3s = [data_cuartiles[2], eq1_cuartiles[2], eq2_cuartiles[2]]
32
33 # Sesgo
34 eq1_sesgo = stats.skew(yj_eq1)
35 eq2_sesgo = stats.skew(yj_eq2)
36 data_sesgo = stats.skew(data)
37 sesgos = [data_sesgo, eq1_sesgo, eq2_sesgo]
38
39 # Curtosis
40 eq1_curt = stats.kurtosis(yj_eq1)
41 eq2_curt = stats.kurtosis(yj_eq2)
42 data_curt = stats.kurtosis(data)
43 curts = [data_curt, eq1_curt, eq2_curt]
44
45 # Tabla
46 tabla_info = {'Mínimo': mins, 'Q1': q1s, 'Mediana': medians, 'Media': means, 'Q3': q3s, 'Máximo': máximos}
47 tabla = pd.DataFrame(tabla_info)
48 print(tabla)

```

	Mínimo	Q1	Mediana	Media	Q3	Máximo
0	1.000000	96.000000	191.000000	4.575686e+02	8.310000e+02	1.801000e+03
1	1.000000	203410.056216	811813.074954	1.083886e+06	1.825210e+06	3.243601e+06
2	0.744581	12.165260	14.750005	1.389156e+01	1.644228e+01	1.773135e+01

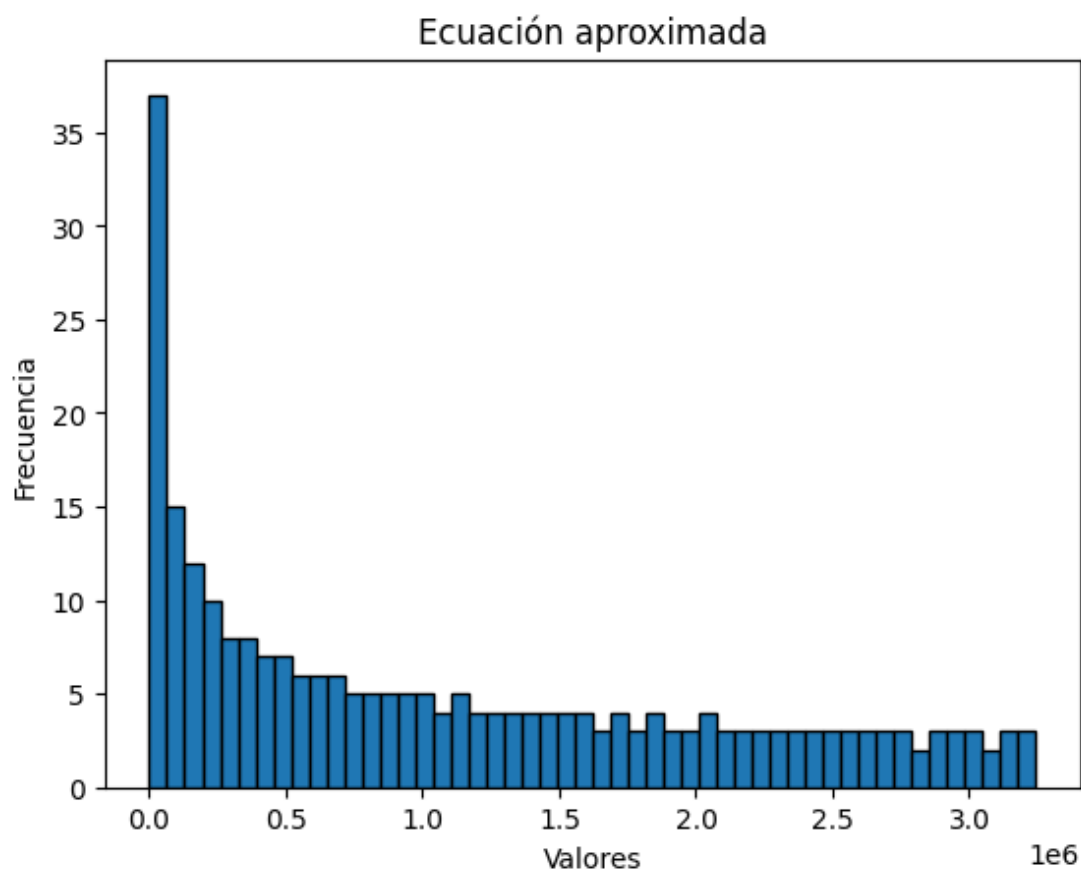
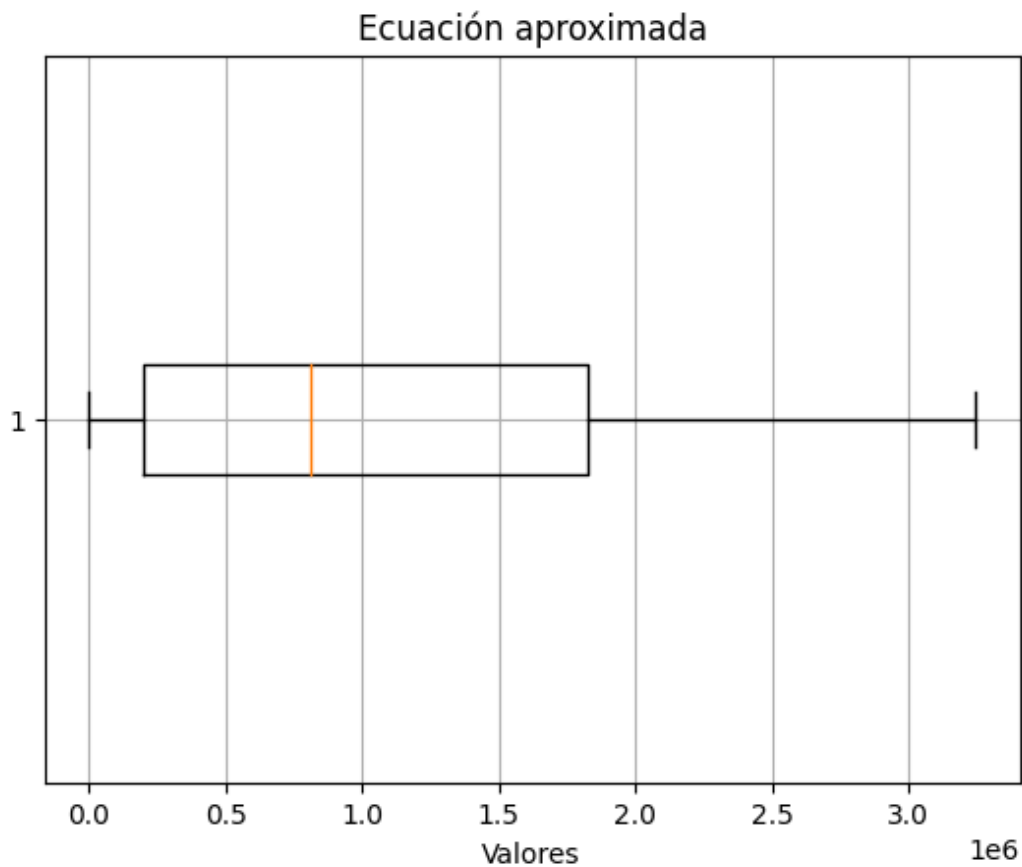
	Curtosis
0	-0.401472
1	-0.855582
2	1.231824

```

1 # Valores Ecuación aproximada
2 plt.boxplot(yj_eq1, vert=False)
3 plt.title("Ecuación aproximada")
4 plt.xlabel("Valores")
5 plt.grid(True)
6 plt.show()
7
8 # Visualizar el histograma
9 plt.hist(yj_eq1, 50, edgecolor='black')
10 plt.title("Ecuación aproximada")

```

```
11 plt.xlabel("Valores")
12 plt.ylabel("Frecuencia")
13 plt.show()
```



```
1 # Valores Ecuación precisa
2 plt.boxplot(yj_eq2, vert=False)
3 plt.title("Ecuación aproximada")
4 plt.xlabel("Valores")
5 plt.grid(True)
6 plt.show()
7 # Visualizar el histograma
8 plt.hist(yj_eq2, 50, edgecolor='black')
9 plt.title("Ecuación aproximada")
10 plt.xlabel("Valores")
11 plt.ylabel("Frecuencia")
12 plt.show()
```

Ecuación aproximada

--	--	--	--	--	--	--	--

Haz doble clic (o ingresa) para editar

```

1 _yj_eq1, critical_values_eq1, significance_levels_eq1 = stats.anderson(yj_eq1, dist='norm')
2 _yj_eq2, critical_values_eq2, significance_levels_eq2 = stats.anderson(yj_eq2, dist='norm')
3 _data, critical_values_data, significance_levels_data = stats.anderson(data, dist='norm')
4
5 print(f"Resultados para Ecuación Aproximada:")
6 print(f"Estadístico Anderson-Darling: {_yj_eq1}")
7 print("Niveles críticos:", critical_values_eq1)
8 print("Niveles de significancia:", significance_levels_eq1)
9 print("-" * 30)
10
11 print(f"\nResultados para Ecuación Exacto:")
12 print(f"Estadístico Anderson-Darling: {_yj_eq2}")
13 print("Niveles críticos:", critical_values_eq2)
14 print("Niveles de significancia:", significance_levels_eq2)
15 print("-" * 30)
16
17 print(f"\nResultados para Datos originales:")
18 print(f"Estadístico Anderson-Darling: {_data}")
19 print("Niveles críticos:", critical_values_data)
20 print("Niveles de significancia:", significance_levels_data)
21 print("-" * 30)

```

➡ Resultados para Ecuación Aproximada:
 Estadístico Anderson-Darling: 8.767656369418944
 Niveles críticos: [0.567 0.646 0.775 0.904 1.076]
 Niveles de significancia: [15. 10. 5. 2.5 1.]

Resultados para Ecuación Exacto:
 Estadístico Anderson-Darling: 7.324450023099871
 Niveles críticos: [0.567 0.646 0.775 0.904 1.076]
 Niveles de significancia: [15. 10. 5. 2.5 1.]

Resultados para Datos originales:
 Estadístico Anderson-Darling: 22.737385535393685
 Niveles críticos: [0.567 0.646 0.775 0.904 1.076]
 Niveles de significancia: [15. 10. 5. 2.5 1.]

Vemos que la diferencia entre el valor estadístico y el valor crítico es muy importante, por lo que no se sigue una distribución normal.

Se tuvieron resultados muy parecidos al modelo de Box-cox.

Conclusiones

El escalamiento permite un mejor ajuste en la distribución de los datos debido a que se permite la modificación de las características de forma de las variables a conveniencia para que se compartan valores medios o variaciones. Forzar la normalización de los datos puede llevarnos a un cambio radical en la distribución de los mismos, obteniendo distribuciones contrarias a la vista en la distribución de los datos originales.

[Productos pagados de Colab](#) - [Cancela los contratos aquí](#)

✓ 0 s se ejecutó 05:26

