

```
1 # Bibliotecas básicas
2 import pandas as pd
3 import numpy as np
4 import matplotlib.pyplot as plt
```

```
1 # Desviación estándar
2 # Conjunto ficticio con distribución normal, notar la 'n' al final
3 np.random.seed(0)
4 data = np.random.randn(50000) * 20 + 20
5 data[1], data[-1]
```

➡ (np.float64(28.003144167344466), np.float64(-5.016539269717178))

```
1 # Anomalías con std VERSIÓN PROPIA
2 corte = 3
3 std = np.std(data)
4 media = np.mean(data)
5 #           antes de 3 std           después de 3 std
6 anom = data[(data < media-corte*std) | (data > media+corte*std)]
7 len(anom)
```

➡ 148

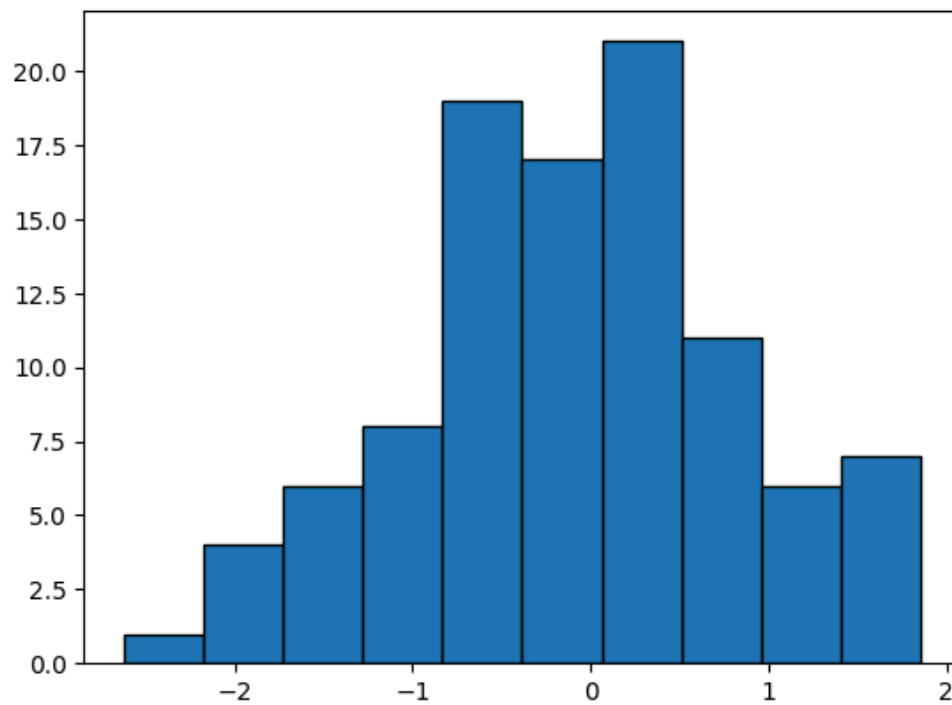
```
1 # Anomalías con std VERSIÓN DE LALO
2 corte = 3
3 std = np.std(data)
4 media = np.mean(data)
5 #           antes de 3 std           después de 3 std
6 anom = data[(data < media-corte*std) | (data > media+corte*std)]
7 len(anom)
```

➡ 148

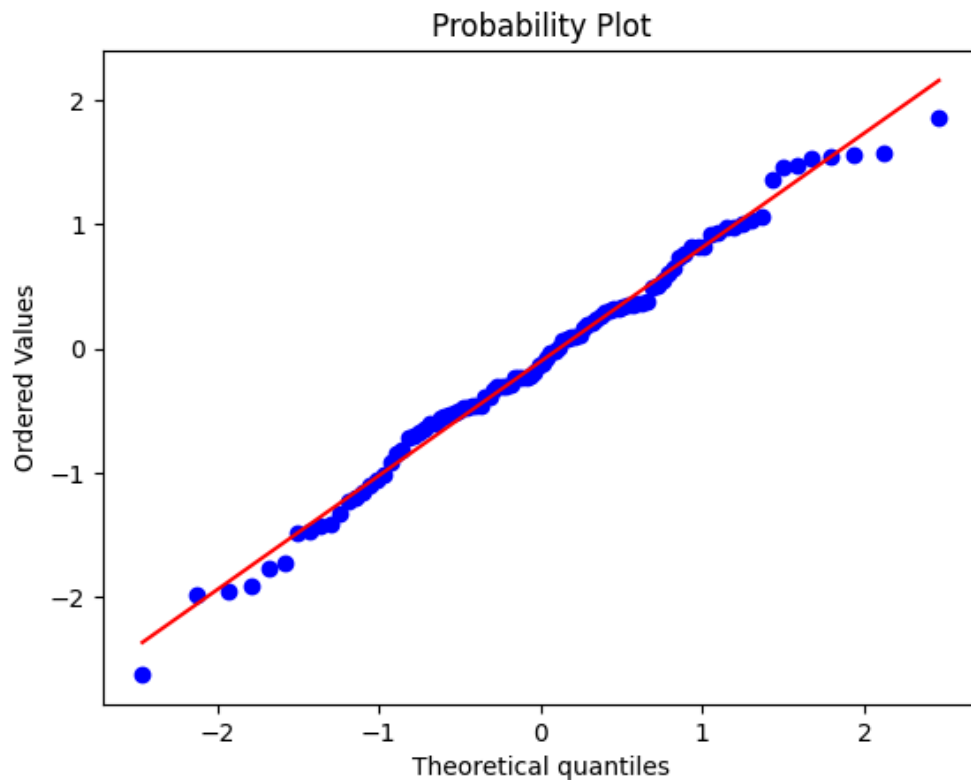
```
1 # Prueba de normalidad: histograma, Quantile-quantile graph, Shapiro
2 # Conjunto de datos ficticio normal
3 np.random.seed(42)
4 data = np.random.normal(0,1,100)
```

```
1 # 1. Histograma
2 plt.hist(data, edgecolor='black', linewidth=1)
```

```
⇒ (array([ 1.,  4.,  6.,  8., 19., 17., 21., 11.,  6.,  7.]),  
   array([-2.6197451, -2.17254278, -1.72534045, -1.27813812, -0.83093579,  
          -0.38373346,  0.06346887,  0.5106712 ,  0.95787353,  1.40507586,  
          1.85227818])),  
   <BarContainer object of 10 artists>)
```



```
1 # 2. gráfica de cuantiles  
2 import pylab  
3 import scipy.stats as stats  
4 prueba = stats.probplot(data, dist='norm', plot=pylab)  
5 pylab.show()
```



```
1 # prueba
```

```
1 # 3. Prueba de Shapiro
2 from scipy.stats import shapiro
3 estad, p_value = shapiro(data)
4 print('p_value = ',p_value)
5 # p_value > 0.05 => distribución normal
```



```
p_value = 0.6551676754214818
```

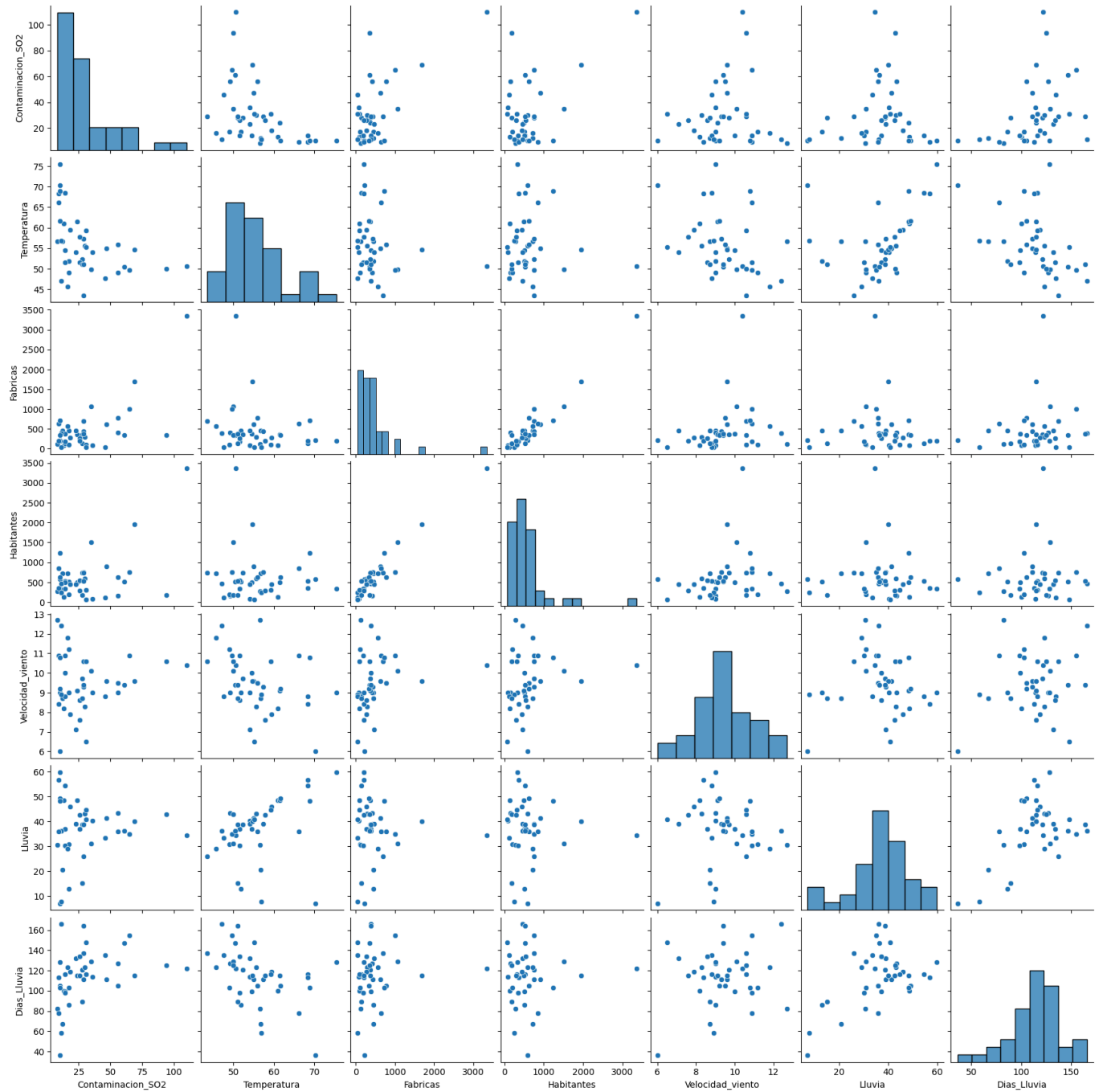
```
1 # Conjunto de datos real
2 cont = pd.read_csv('https://bit.ly/31B56KB')
3 cont.info()
```



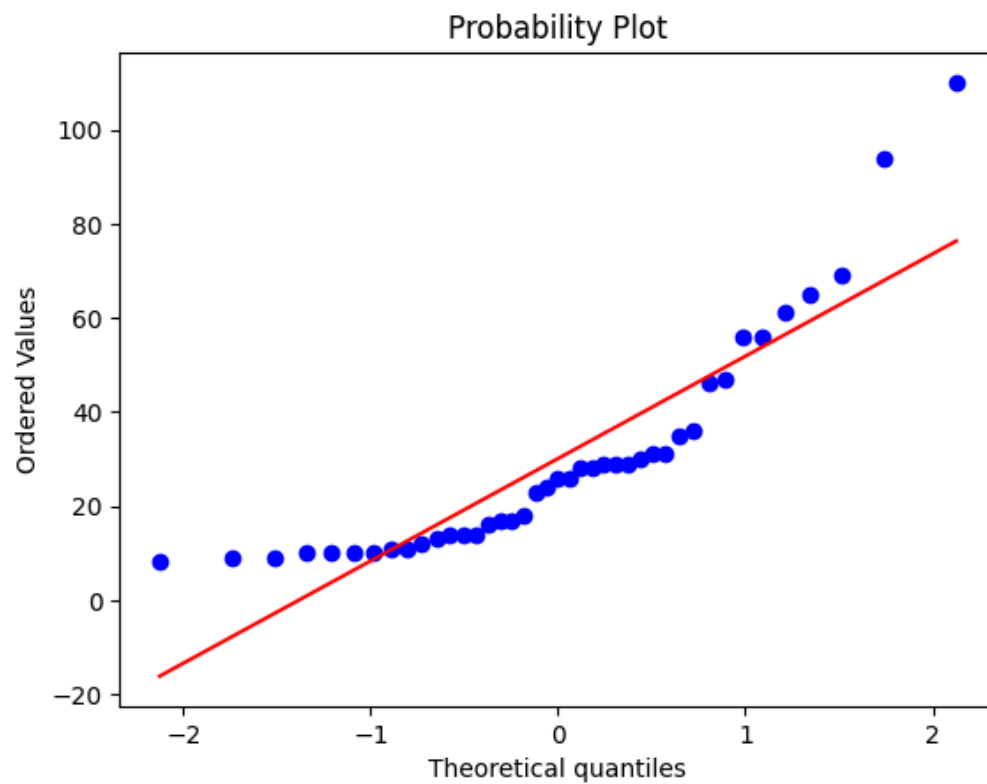
```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 41 entries, 0 to 40
Data columns (total 7 columns):
#   Column                Non-Null Count  Dtype
---  -
0   Contaminacion_S02      41 non-null    int64
1   Temperatura            41 non-null    float64
2   Fabricas               41 non-null    int64
3   Habitantes             41 non-null    int64
4   Velocidad_viento       41 non-null    float64
5   Lluvia                 41 non-null    float64
6   Dias_Lluvia            41 non-null    int64
dtypes: float64(3), int64(4)
memory usage: 2.4 KB
```

```
1 import seaborn as sns
2 sns.pairplot(cont)
```

↔ <seaborn.axisgrid.PairGrid at 0x7f0ec084e590>



```
1 # Contaminación_SO2
2 # Quartiles
3 stats.probplot(cont.Contaminacion_SO2, dist='norm', plot=pylab)
4 pylab.show()
```

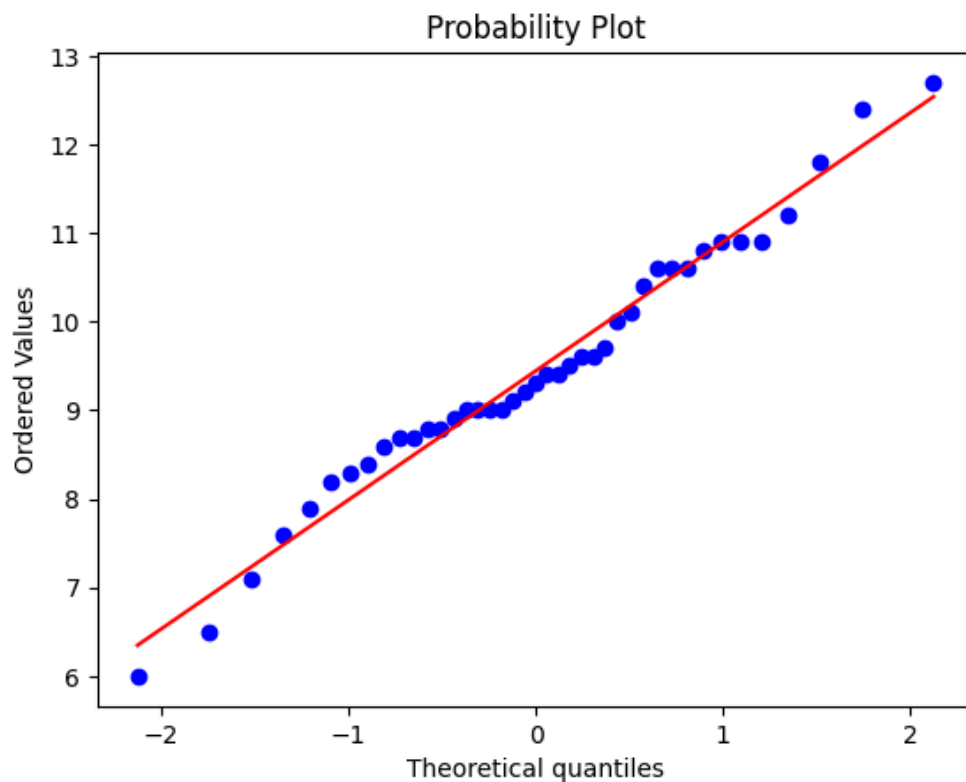


```
1 # Shapiro
2 estad, p_value = shapiro(cont.Contaminacion_S02)
3 print('p_value = ', p_value)
```



```
p_value = 9.723376400158156e-06
```

```
1 # Velocidad_viento
2 stats.probplot(cont.Velocidad_viento, dist='norm', plot=pylab)
3 pylab.show()
```



```
1 # Shapiro
2 estad, p_value = shapiro(cont.Velocidad_viento)
3 print('p_value = ', p_value)
```



p_value = 0.6972579783041465

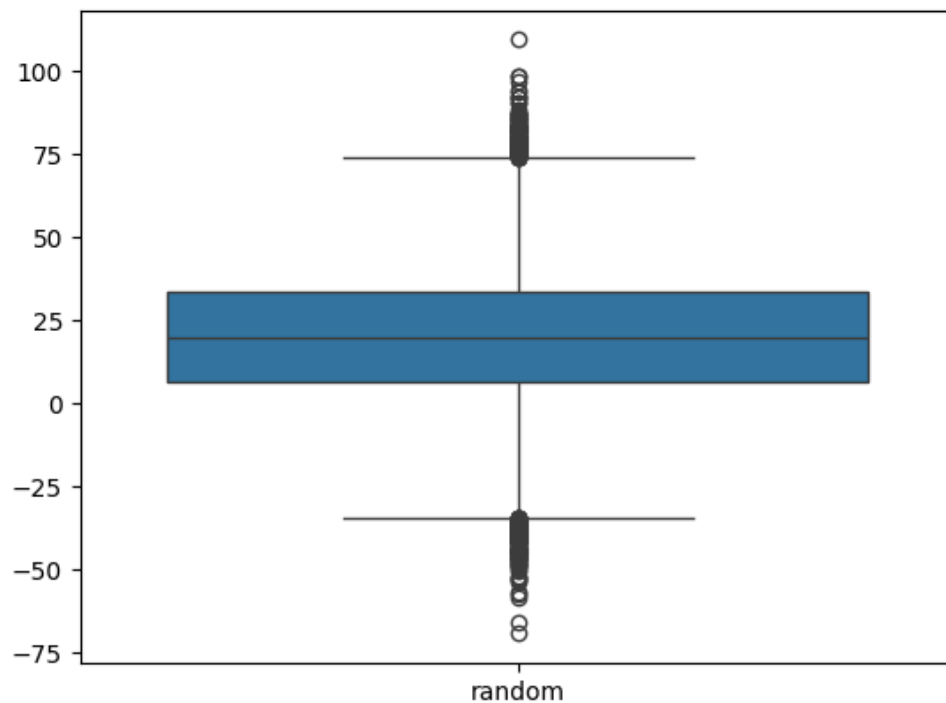
```
1 # Diagrama de bigotes para datos con distribución normal
2 data = np.random.randn(50000) * 20 + 20
3 data = pd.DataFrame(data, columns=['random'])
4 data.head(2)
```



	random
0	-8.307415
1	11.587094

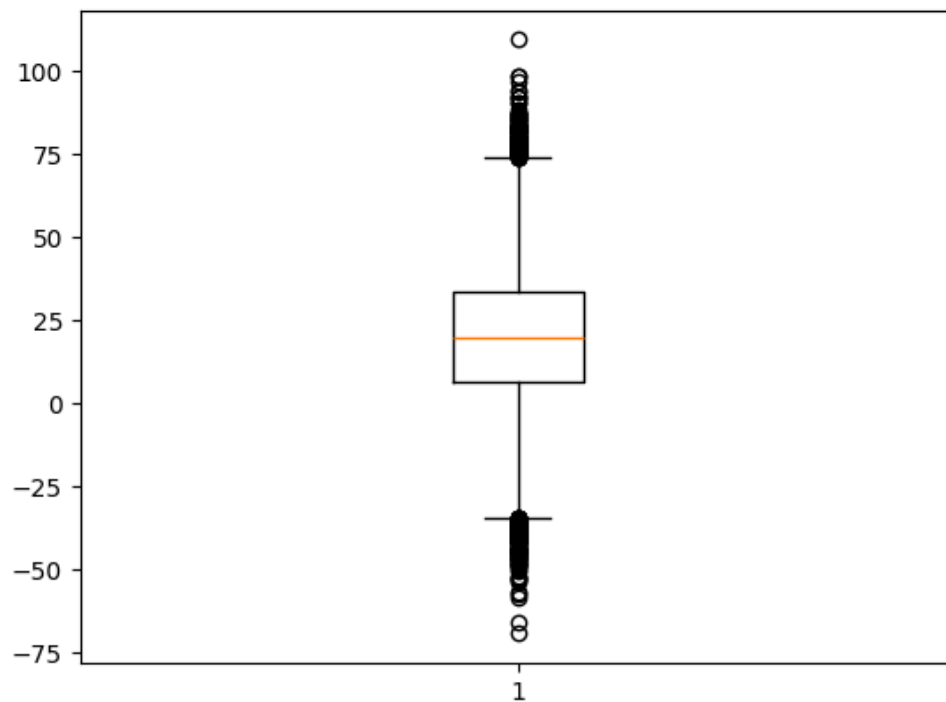
```
1 import seaborn as sns
2 sns.boxplot(data)
```

↔ <Axes: >



```
1 bp = plt.boxplot(x=data)
```

↔



1 bp

↔

```
{'whiskers': [<matplotlib.lines.Line2D at 0x7f0ebdf7dad0>,  
<matplotlib.lines.Line2D at 0x7f0ebc28a410>],  
'caps': [<matplotlib.lines.Line2D at 0x7f0ebc23ff10>,  
<matplotlib.lines.Line2D at 0x7f0ebc23db90>],  
'boxes': [<matplotlib.lines.Line2D at 0x7f0ebdf7e550>],  
'medians': [<matplotlib.lines.Line2D at 0x7f0ebc23dcd0>],
```

```
'fliers': [<matplotlib.lines.Line2D at 0x7f0ebc23cfd0>],
'means': []}
```

```
1 bp['fliers']
```

```
→ [<matplotlib.lines.Line2D at 0x7f0ebc23cfd0>]
```

```
1 #bp['fliers'][0].get_data()[1]
```

```
1 bp['fliers'][0].get_data()[1].shape
```

```
→ (332,)
```

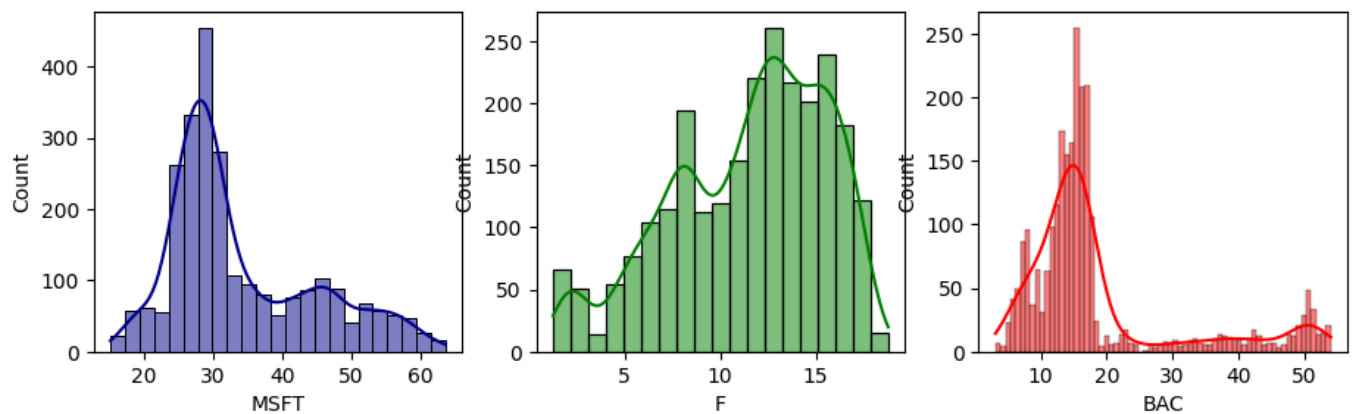
```
1 stocks = pd.read_csv('stocks.csv', header='infer')
2 stocks.index = stocks.Date
3 #stocks.set_index(stocks.Date)
4 stocks.drop(['Date'],axis=1, inplace=True)
5 stocks.head(3)
```

```
→
```

	MSFT	F	BAC
Date			
1/3/2007	29.860001	7.51	53.330002
1/4/2007	29.809999	7.70	53.669998
1/5/2007	29.639999	7.62	53.240002

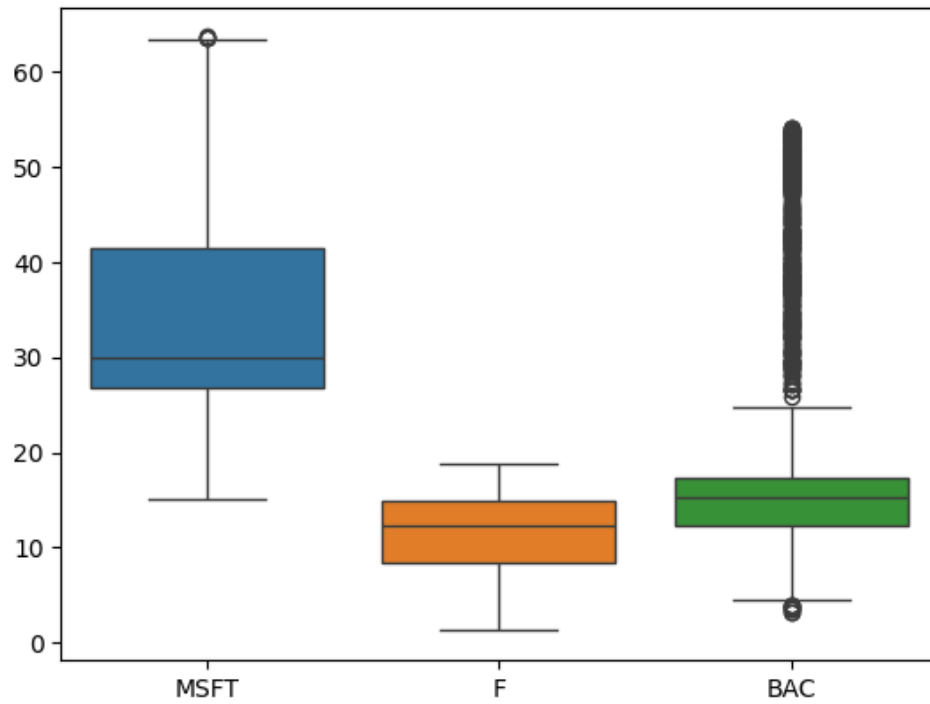
```
1 import matplotlib.pyplot as plt
2 import seaborn as sns
3 fig, ax = plt.subplots(nrows=1, ncols=3, figsize=(11,3))
4 sns.histplot(stocks.MSFT, ax=ax[0], color="darkblue", kde=True)
5 sns.histplot(stocks.F, ax=ax[1], color='green', kde=True)
6 sns.histplot(stocks.BAC, ax=ax[2], color='red', kde=True)
```

```
→ <Axes: xlabel='BAC', ylabel='Count'>
```



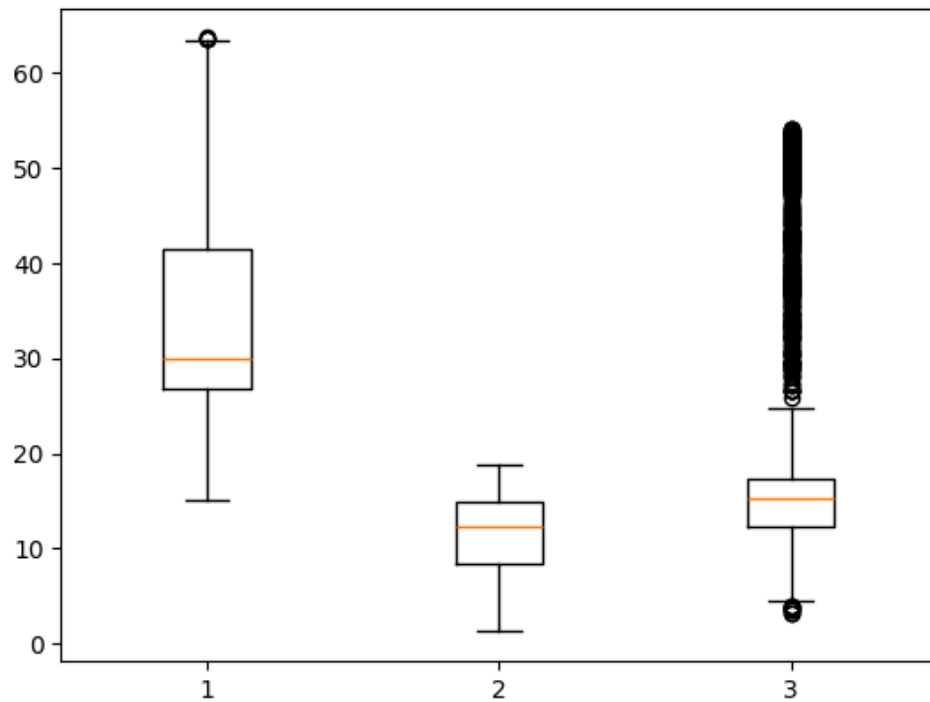
```
1 sns.boxplot(data=stocks)
```


 <Axes: >



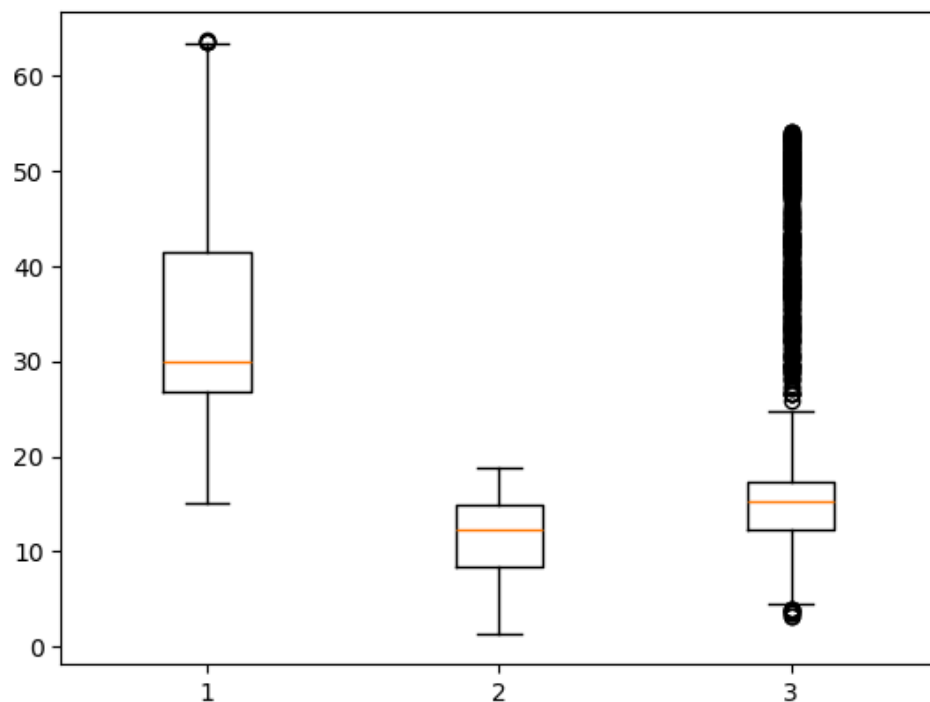
```
1 bp = plt.boxplot(x=stocks)
```





```
1 bp['fliers'][0].get_data()[1].shape
```

```
1 bp = plt.boxplot(x=stocks)
```



```
1 bp['fliers'][0].get_data()[1].shape
```



```
(4,)
```

```
1 bp['fliers'][1].get_data()[1].shape
```



```
(0,)
```

```
1 bp['fliers'][2].get_data()[1].shape
```



```
(441,)
```

```
1
```

```
1 # Obtener los atípicos con IQR
2 np.random.seed(0)
3 data = np.random.randn(50000) * 20 + 20
```

```
1 Q1, Q3 = np.percentile(data, [25, 75])
2 IQR = Q3-Q1
3 # bigotes (whiskers)
4 sup = Q3 + 1.5*IQR
5 inf = Q1 - 1.5*IQR
6 # anomalías
7 anom_ind = np.where((data < inf) | (data > sup))
8 anom = data[anom_ind] # Acceso directo a las anomalías
9 len(anom_ind[0]), len(anom)
```



```
(379, 379)
```

```

1
1 # EllipticEnvelope:
2 # Basado en covarianza, genera una elipse que "envuelve" los datos "normales"
3 # y deja fuera los típicos
4 import numpy as np
5 import matplotlib.pyplot as plt
6 from sklearn.covariance import EllipticEnvelope
7 cov = np.array([[.8, .3],
8                 [.3, .4]])
9 X = np.random.RandomState(0).multivariate_normal(mean=[0,0],cov=cov,size=500)
10 ee = EllipticEnvelope(random_state=42) # random_state parámetro de replicabilidad
11 ee.fit(X)
12 # predict regresa 1 para valores normales y -1 para atípicos
13 anom = ee.predict(X)
14 ee.covariance_, ee.location_

```

```

➡ (array([[0.7361743 , 0.25159568],
          [0.25159568, 0.30556494]]),
   array([0.07405583, 0.03661301]))

```

```

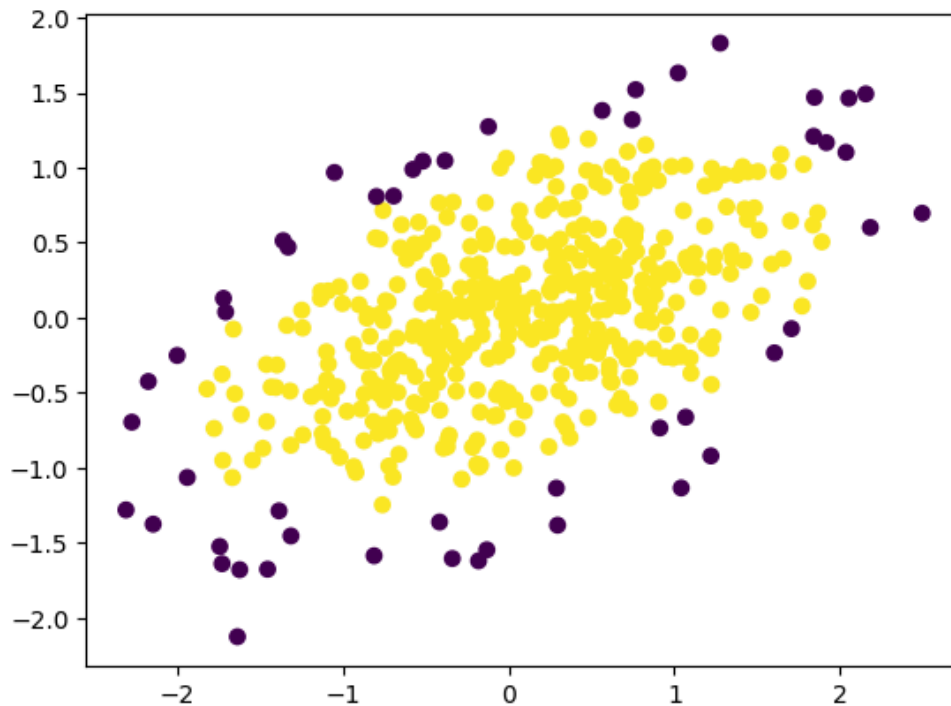
1 plt.scatter(X[:,0],X[:,1],c=anom)

```

```

➡ <matplotlib.collections.PathCollection at 0x7f0eb4471310>

```



```

1 #anom

```

```

1 # Funciona para datos con múltiples dimensiones

```

```

1

```