

Outlier Detection

Aim: To detect outliers in a dataset using the Interquartile Range (IQR) method and visualize them using a scatter plot.

Algorithm:

Outlier detection refers to the process of identifying data points that significantly differ from the rest of the dataset.

One common method of outlier detection is using the Interquartile Range (IQR).

Calculate the Interquartile Range (IQR): The IQR is the difference between the third quartile (Q3) and the first quartile (Q1). It measures the spread of the middle 50% of the data.

$$IQR = Q3 - Q1$$

Determine the lower and upper bounds for outliers: Outliers are defined as data points that lie below the lower bound or above the upper bound. These bounds are calculated as:

- **Lower Bound:**

$$\text{Lower Bound} = Q1 - 1.5 \times IQR$$

- **Upper Bound:**

$$\text{Upper Bound} = Q3 + 1.5 \times IQR$$

Identify outliers: Any data point that is less than the lower bound or greater than the upper bound is considered an outlier. These points are usually isolated and significantly different from the majority of the dataset.

Step 1: Import Libraries

- Import necessary Python libraries such as pandas for data manipulation, numpy for numerical operations, and matplotlib for visualization.

Step 2: Generate the Dataset

- Create a sample dataset using `numpy.random.randn()` to generate random data points and introduce outliers using `numpy.random.uniform()`.

Step 3: Prepare the Data

- Convert the data into a pandas DataFrame for easier manipulation and visualization. This allows for better handling of outliers and feature selection.

Step 4: Calculate the Quartiles

- Calculate the 25th percentile (Q1) and the 75th percentile (Q3) of the data using the `quantile()` function of pandas.

Step 5: Calculate the Interquartile Range (IQR)

- Calculate the IQR as the difference between Q3 and Q1:

$$IQR = Q3 - Q1$$

Step 6: Determine Lower and Upper Bounds for Outliers

- Calculate the lower and upper bounds using the IQR:
- Lower Bound:

$$\text{Lower Bound} = Q1 - 1.5 \times IQR$$

- Upper Bound:

$$\text{Upper Bound} = Q3 + 1.5 \times IQR$$

Step 7: Identify the Outliers

- Identify data points that fall outside of the lower and upper bounds as outliers. These are points that are significantly different from the majority of the dataset.

Step 8: Visualize the Data

- Plot the data points using `matplotlib` where inliers are shown in one color and outliers are shown in another color, visually identifying the outliers in the dataset.

Import necessary libraries

```
In [24]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
```

Generate a sample dataset

```
In [25]: np.random.seed(42)
data = np.random.randn(700, 2) * 10 + 50

In [26]: outliers = np.random.uniform(low=20, high=80, size=(10, 2)) # 10 random outliers
data_with_outliers = np.vstack([data, outliers])
```

Convert data to DataFrame

```
In [27]: df = pd.DataFrame(data_with_outliers, columns=['Feature1', 'Feature2'])
```

```
In [28]: df
```

Out[28]:

	Feature1	Feature2
0	54.967142	48.617357
1	56.476885	65.230299
2	47.658466	47.658630
3	65.792128	57.674347
4	45.305256	55.425600
...
705	37.726687	66.153389
706	57.479814	42.916378
707	32.341236	27.283185
708	56.900778	66.478027
709	58.634255	51.818128

710 rows × 2 columns

Calculate 25th and 75th percentile

```
In [29]: Q1 = df.quantile(0.25)
Q3 = df.quantile(0.75)
```

Calculate Interquartile Range

```
In [30]: IQR = Q3 - Q1
```

Set lower and upper bounds for outliers

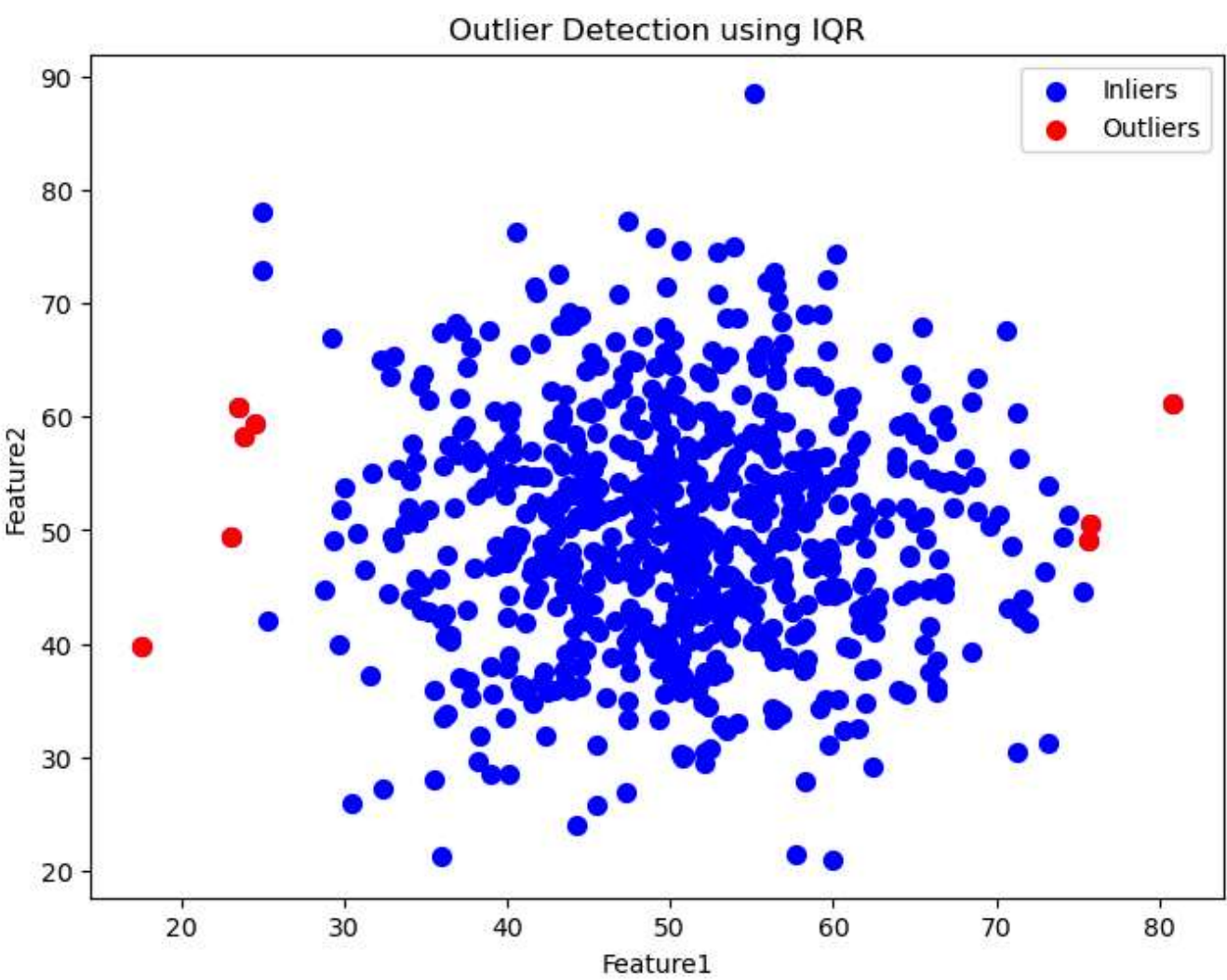
```
In [31]: lower_bound = Q1 - 1.5 * IQR
upper_bound = Q3 + 1.5 * IQR
```

Identify the outliers

```
In [32]: outliers = (df < lower_bound) | (df > upper_bound)
```

Plot the graph

```
In [33]: plt.figure(figsize=(8, 6))
plt.scatter(df['Feature1'], df['Feature2'], color='blue', label='Inliers', s=50)
plt.scatter(df[outliers['Feature1']]['Feature1'], df[outliers['Feature1']]['Feature2'], color='red', label='Outliers', s=50)
plt.title('Outlier Detection using IQR')
plt.xlabel('Feature1')
plt.ylabel('Feature2')
plt.legend()
plt.show()
```



```
In [34]: outliers_data = df[outliers.any(axis=1)]
print("Outliers detected:")
print(outliers_data)
```

Outliers detected:

	Feature1	Feature2
37	23.802549	58.219025
104	55.150477	88.527315
131	17.587327	39.756124
239	80.788808	61.195749
323	23.031134	49.457051
327	75.733598	50.592184
334	23.490302	60.915069
381	75.600845	49.039401
530	57.716987	21.514574
550	59.980101	21.037446
673	24.460789	59.343199
703	35.902786	21.256698
704	24.930300	78.071602

Result

Outlier detection was successfully performed on a dataset using the Interquartile Range (IQR) method. The identified outliers were isolated and visualized, highlighting 13 points as outliers.