

# IQR technique to handle the outliers

## Example 1

Here's a Python code using the numpy library to demonstrate the IQR method for handling outliers:

```
In [86]: import numpy as np

def detect_outliers_iqr(data, factor=1.5):
    """
    Detects outliers using IQR method.

    :param data: List or array-like data points.
    :param factor: Multiplier for IQR, usually 1.5. Adjust to be more or less strict.
    :return: List of outliers.
    """
    quartile_1, quartile_3 = np.percentile(data, [25, 75])
    iqr = quartile_3 - quartile_1
    lower_bound = quartile_1 - (iqr * factor)
    upper_bound = quartile_3 + (iqr * factor)

    outliers = [point for point in data if point < lower_bound or point > upper_bound]
    return outliers

# data points
data = [10, 12, 12, 13, 12, 11, 11, 52, 13, 12, 11]

# Detecting outliers
outliers = detect_outliers_iqr(data)

print("Outliers:", outliers)

Outliers: [52]
```

In the above code:

We first compute the first and third quartiles. We then compute the IQR. We identify the lower and upper bounds for what's considered a non-outlier. We finally return the list of values that are below the lower bound or above the upper bound as outliers.

## Example 2

```
In [90]: import pandas as pd
import seaborn as sns

# Load the diamonds dataset from seaborn
df = sns.load_dataset('diamonds')

# Calculate Q1, Q3, and IQR for the 'price' column
Q1 = df['price'].quantile(0.25)
Q3 = df['price'].quantile(0.75)
IQR = Q3 - Q1

# Define bounds
lower_bound = Q1 - 1.5 * IQR
upper_bound = Q3 + 1.5 * IQR

# Filter rows in dataframe to exclude data points that are outliers for 'price'
```

```
df_no_outliers = df[(df['price'] >= lower_bound) & (df['price'] <= upper_bound)]

print("Original Dataset Length:", len(df))
print("\nDataset Length without Outliers based on price:", len(df_no_outliers))
```

Original Dataset Length: 53940

Dataset Length without Outliers based on price: 50400

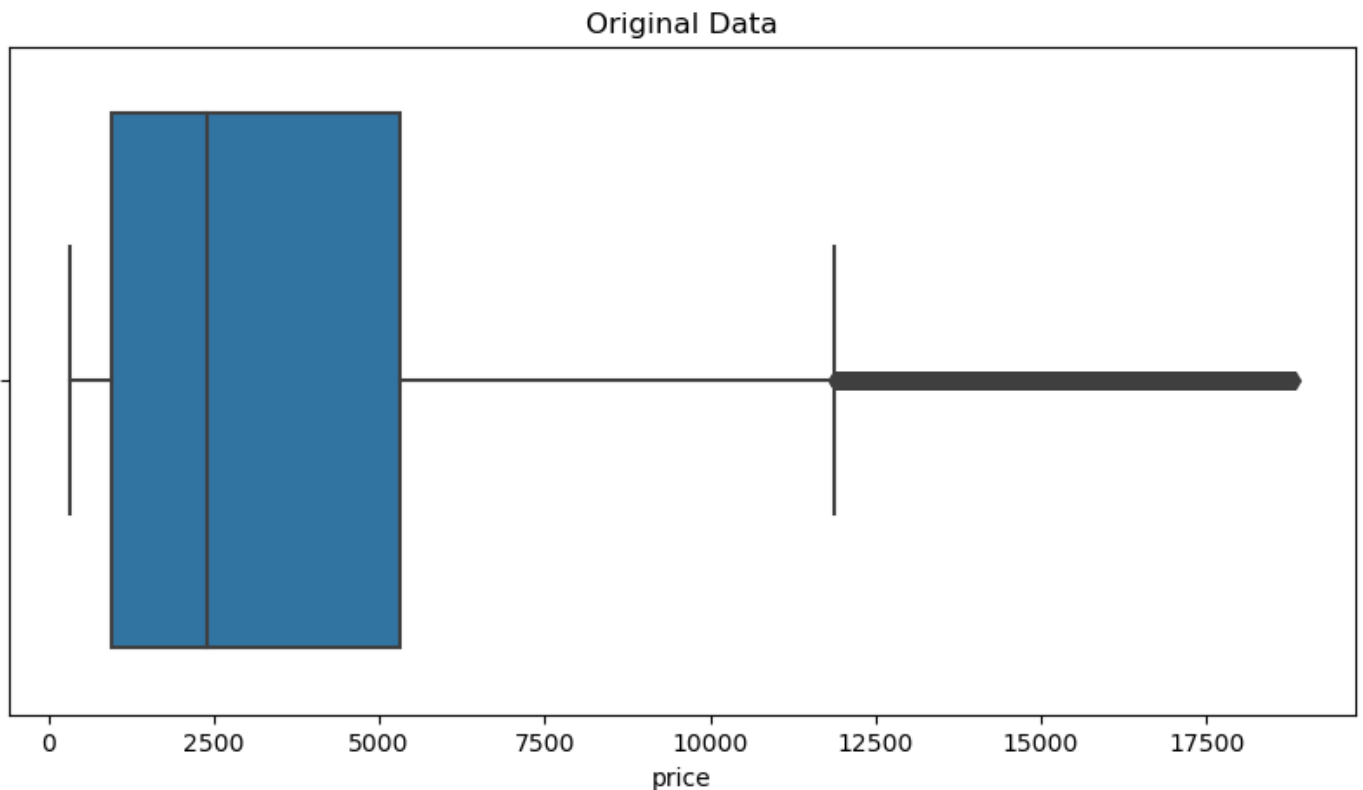
The output will show a difference in the lengths of the original dataset and the filtered dataset, indicating that some outlier rows based on the price column were removed.

You can also visualize the outliers using box plots. Before and after the filtering process:

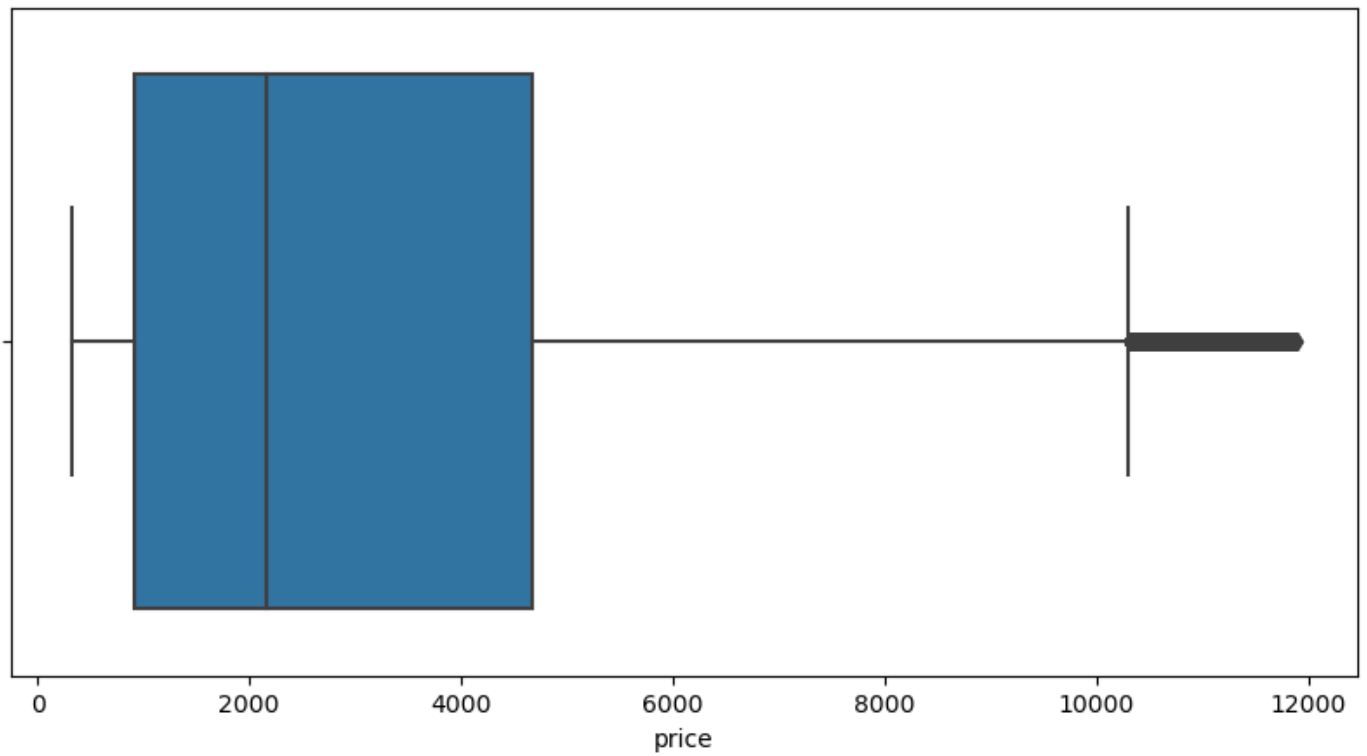
```
In [93]: import matplotlib.pyplot as plt

# Original data visualization
plt.figure(figsize=(10, 5))
sns.boxplot(x=df['price'])
plt.title('Original Data')
plt.show()

# Data without outliers visualization
plt.figure(figsize=(10, 5))
sns.boxplot(x=df_no_outliers['price'])
plt.title('Data without Outliers')
plt.show()
```



Data without Outliers



```
In [95]: # import sklearn module
import sklearn
from sklearn import datasets

# Loading california dataset for getting outliers using IQR method
california_dateset = sklearn.datasets.fetch_california_housing(as_frame=True)
df = pd.DataFrame(california_dateset.data)
```

```
In [97]: df.head(5)
```

```
Out[97]:
```

	MedInc	HouseAge	AveRooms	AveBedrms	Population	AveOccup	Latitude	Longitude
0	8.3252	41.0	6.984127	1.023810	322.0	2.555556	37.88	-122.23
1	8.3014	21.0	6.238137	0.971880	2401.0	2.109842	37.86	-122.22
2	7.2574	52.0	8.288136	1.073446	496.0	2.802260	37.85	-122.24
3	5.6431	52.0	5.817352	1.073059	558.0	2.547945	37.85	-122.25
4	3.8462	52.0	6.281853	1.081081	565.0	2.181467	37.85	-122.25

```
In [99]: df.info()
```



```
20428    -118.83
20436    -118.69
20462    -118.76
```

```
[511 rows x 8 columns]
```

```
In [81]: !jupyter nbconvert --to webpdf --allow-chromium-download Week8_Lab2.ipynb
```

```
[NbConvertApp] Converting notebook Week8_Lab2.ipynb to webpdf
```

```
[NbConvertApp] Building PDF
```

```
[NbConvertApp] PDF successfully created
```

```
[NbConvertApp] Writing 342267 bytes to Week8_Lab2.pdf
```

```
In [ ]:
```