# Handling Outliers using Python:

**Colab:**

Outliers are data points that significantly differ from other observations in a dataset. They can affect statistical analysis and machine learning models, so handling them properly is important.

**Common methods to handle outliers in Python:**

1. **Detection using IQR (Interquartile Range)**

   o Calculate Q1, Q3, and IQR.

   o Define lower and upper fences using Q1 - 1.5×IQR and Q3 + 1.5×IQR.

   o Values outside these limits are treated as outliers.

2. **Detection using Z-Score**

   o Measures how many standard deviations a value is from the mean.

   o Typically, values with a Z-score greater than ±3 are considered outliers.

3. **Removing Outliers**

   o Outliers can be removed if they are due to errors or noise.

   o This is done by filtering data within acceptable limits.

4. **Capping (Winsorization)**

   o Outliers are replaced with the nearest boundary values instead of removing them.

   o Useful when data loss is undesirable.

5. **Transformation**

   o Applying transformations like logarithmic or square root can reduce the impact of outliers.

**Short Notes on IQR (Interquartile Range)**

**Definition:**
The Interquartile Range (IQR) is the range of the middle 50% of a dataset. It measures statistical dispersion and is calculated as:

$$\text{IQR} = Q3 - Q1$$

Where:

- **Q1** = First Quartile (25th percentile)

- **Q3** = Third Quartile (75th percentile)

**Purpose:**

- Measures the spread of the central portion of the data.

- Less sensitive to extreme values compared to range or standard deviation.

- Used for **outlier detection** using the formula:

    - Lower Fence = Q1 − 1.5 × IQR

    - Upper Fence = Q3 + 1.5 × IQR