# Comparative Report: IQR vs. Z-Score for Outlier Detection

## Introduction

Outlier detection is crucial in data analysis to ensure the integrity and accuracy of insights derived from datasets. Two popular methods for detecting outliers are the Interquartile Range (IQR) and Z-Score. This report compares these two methods in terms of their strengths and weaknesses.

## IQR Method

**Pros:**

1. **Robustness to Non-Normal Distributions**:
   - IQR is based on percentiles and is less affected by the distribution shape, making it suitable for skewed data.
2. **Simple Interpretation**:
   - The concept of lower and upper bounds (Q1 - 1.5 * IQR and Q3 + 1.5 * IQR) is straightforward and easy to understand.
3. **No Assumption of Normality**:
   - IQR does not require data to follow a normal distribution, making it applicable to a wider variety of datasets.

**Cons:**

1. **Sensitivity to Sample Size**:
   - IQR may not perform well with small sample sizes, where the estimate of the quartiles can be unstable.
2. **Limited to Linear Relationships**:
   - IQR can miss outliers in more complex distributions where the relationship is not linear.
3. **Static Threshold**:
   - The standard multiplier (1.5) may not be optimal for all datasets and could require manual adjustments.

## Z-Score Method

**Pros:**

1. **Sensitivity to Standard Deviation**:
   - Z-Score effectively highlights outliers based on how many standard deviations away a data point is from the mean, providing a relative perspective.
2. **Dynamic Thresholds**:
   - By allowing different thresholds (e.g., 2, 2.5, 3), it offers flexibility to tune the sensitivity of outlier detection based on specific dataset characteristics.
3. **Works Well with Normal Distributions**:
   - It is particularly effective when data is normally distributed, providing a clear mechanism for detecting deviations.

**Cons:**

1. **Assumption of Normality**:
   - Z-Scores assume that the data follows a normal distribution; this can lead to misleading results if the assumption is violated.
2. **Sensitivity to Outliers**:
   - The presence of outliers can skew the mean and standard deviation, which may distort the Z-Score calculation itself.
3. **Not Robust to Non-Normal Distributions**:
   - It can misidentify outliers in datasets that are highly skewed or have heavy tails.

## Comparative Analysis

The comparison of anomalies detected by both methods across multiple parameters highlights their performance characteristics. Here's a summary of key findings:

1. **Detection Variability**:
   - For certain parameters, the IQR method may identify outliers that the Z-Score method does not, especially in non-normally distributed data. Conversely, Z-Score may detect subtle anomalies in normally distributed parameters.

2. **Scalability**:
   - o Scaling the data using Min-Max scaling before applying Z-Scores allows for a more standardized analysis, though it introduces additional preprocessing steps.
3. **Visualization**:
   - o Visual tools like scatter plots provide a valuable means to understand the performance of both methods visually, showing overlaps and discrepancies in detected anomalies.

## Conclusion

The choice between IQR and Z-Score for outlier detection largely depends on the characteristics of the dataset and the specific analysis requirements.

- **IQR** is preferred for robustness against **non-normality and simplicity**.
- **Z-Score** is advantageous when the data is approximately **normally** distributed and when **dynamic thresholds** are required.