## Computer Networks: Detecting and Treating Outliers

**I. Introduction**

* **Definition:** An outlier is a data point that significantly deviates from other data points in a dataset.
* **Causes:** Outliers can arise from various factors:
* Measurement errors
* Data entry errors
* Natural variations within the data
* **Impact:** Outliers can significantly influence:
* Machine learning analysis and model performance.
* Data interpretation and the drawing of conclusions.
* **Importance:** Detecting and handling outliers is crucial for accurate data analysis and reliable results.

**II. Outlier Detection Methods**

**A. Standard Deviation Method**

* **Concept:** Identify outliers by measuring how far data points are from the mean, using standard deviations.
* **Process:**
1. Calculate the standard deviation ($\sigma$) of the dataset.
2. Calculate the mean ($\mu$) of the dataset.
3. Identify data points that lie outside a specified number of standard deviations (e.g., 2 or 3) from the mean ($\mu \pm k\sigma$).
* **Limitations:** This method assumes the data follows a normal distribution. It might not be suitable for non-normally distributed datasets.

**B. Z-Score Method**

* **Concept:** The Z-score measures how many standard deviations a data point is away from the mean.
* **Formula:**
* `z = (x - μ) / σ`
* `z`: Z-score
* `x`: Value of the data point
* `μ`: Mean of the dataset
* `σ`: Standard deviation of the dataset
* **Z-Score Interpretation:**
* `z = 0`: The data point (x) is equal to the mean ($\mu$).
* `z = ±1`: The data point is one standard deviation away from the mean.
* `z = ±2`: The data point is two standard deviations away from the mean.
* `z = ±3`: The data point is three standard deviations away from the mean.
* **Outlier Identification:** A data point is typically considered an outlier if its Z-score is greater than 3 or less than -3.
* **Assumptions:** This method assumes the data follows a normal distribution.

**C. Interquartile Range (IQR) Method**

* **Concept:** Identifies outliers based on the data's quartiles and the IQR.
* **Steps:**

1. **Sort the dataset:** Arrange the data in ascending order.

2. **Calculate Quartiles:**

* Q1 (First Quartile): The 25th percentile (median of the lower half of the data).

* Q3 (Third Quartile): The 75th percentile (median of the upper half of the data).

3. **Compute IQR:**

* `IQR = Q3 - Q1`

4. **Calculate Bounds:**

* Lower Bound: `Q1 - 1.5 * IQR`

* Upper Bound: `Q3 + 1.5 * IQR`

5. **Identify Outliers:** Any data points that fall below the lower bound or above the upper bound are considered outliers.

* **Advantages:** Robust to extreme values (less sensitive than methods using mean and standard deviation).

* **Visual Representation:** Box and whisker plots are commonly used to visualize the IQR and identify outliers.

* The box represents the IQR (from Q1 to Q3).

* The whiskers extend to the lowest and highest values within the bounds (or to a certain multiple of the IQR).

* Outliers are plotted as individual points beyond the whiskers.

**III. Example and Application**

**Problem:** Find the outliers, if any, for the following dataset, and draw the box and whisker plot.

**Dataset:** 250, 270, 280, 370, 572, 574, 578, 585, 592, 592, 607, 616, 618, 621, 629, 637, 638, 640, 656, 668, 707, 709, 719, 737, 739, 752, 758, 766, 792, 792, 794, 802, 818, 830, 832, 843, 858, 860, 869, 918, 925, 953, 991, 1000, 1005, 1068, 1441.

**Solution (using IQR method):**

1. **Sorted Dataset:** (already sorted)

2. **Calculate Quartiles:**

* Total Number of data points: 47

* Q1 = Value at position (47 + 1)/4 = 12 -> 616

* Q3 = Value at position 3(47 + 1)/4 = 36 -> 843

3. **Compute IQR:**

* IQR = Q3 - Q1 = 843 - 616 = 227

4. **Calculate Bounds:**

* Lower Bound: Q1 - 1.5 * IQR = 616 - (1.5 * 227) = 616 - 340.5 = 275.5

* Upper Bound: Q3 + 1.5 * IQR = 843 + (1.5 * 227) = 843 + 340.5 = 1183.5

5. **Identify Outliers:**

* Check for values below Lower Bound: 250, 270, 280, 370 are outliers

* Check for values above Upper Bound: 1441 is an outlier

* **Outliers:** 250, 270, 280, 370, 1441

**Box and Whisker Plot (Conceptual):**

(This is a description, imagine a graphical representation)

* **Box:** Drawn from 616 (Q1) to 843 (Q3).

* **Median:** Located within the box (approximately around 739).

* **Whiskers:**

* Left whisker extends from 616 to the lowest data point that is not an outlier (572).

* Right whisker extends from 843 to the highest data point that is not an outlier (1068).

* **Outliers:** Individual points at: 250, 270, 280, 370, 1441 (plotted as dots beyond the whiskers).