

The Importance of Outlier Detection With Applications in R

Rajesh Mahato - 400878749

Fresenius University of Applied Science

Data Analysis for Decision-Making (WS 2025/26)

Prof. Dr. Stephan Huber

2025-12-03

Author Note

Correspondence concerning this article should be addressed to Rajesh Mahato - 400878749, Email: mahato.rajesh@stud.hs-fresenius.de

Abstract

Outlier detection plays a critical role in safeguarding the validity of data-driven decisions by identifying observations that can distort statistical summaries, bias models, or signal rare but important events. This handout focuses on three widely used approaches to outlier detection—the Interquartile Range (IQR) method, Z-score method, and Local Outlier Factor (LOF)—and discusses when each is most appropriate in practice. Using the classical New York air quality dataset as an example, IQR-based detection is implemented in R to flag extreme ozone and solar radiation values, which are then related to public health risk thresholds and temporal pollution patterns. The analysis highlights how univariate methods can be used to build intuitive visual tools (boxplots, time series with highlighted events), while also motivating the need for multivariate techniques such as LOF in complex datasets. Through comparative analysis and practical coding examples, this handout demonstrates that effective outlier detection combines statistical rigor with domain expertise to inform environmental monitoring, fraud prevention, quality control, and predictive modeling across sectors.

The Importance of Outlier Detection With Applications in R

1 Introduction

Outlier detection is a fundamental component of data analysis that has gained substantial relevance across diverse scientific, environmental, and business domains. Outliers—observations that deviate significantly from the expected pattern in a dataset—can profoundly impact statistical analyses, predictive modeling, and downstream business decisions Tukey (1977). This handout expands upon the presentation by providing deeper theoretical foundations, methodological considerations, and practical applications of outlier detection techniques using the R programming language.

2 What Are Outliers?

Outliers are data points that differ significantly from other observations in a dataset. They manifest in various forms and originate from multiple sources, each requiring different handling strategies and interpretations.

2.1 Types of Outliers

- **Point Outliers:** Single extreme values that stand out from the overall distribution.
- **Contextual Outliers:** Observations that are unusual within specific contexts or conditions.
- **Collective Outliers:** Groups of unusual points that, while individually not extreme, form anomalous patterns when considered together.

2.2 Common Causes

- **Measurement errors:** Instrument calibration issues or sensor failures.
- **Data entry mistakes:** Human errors in recording or transcribing data.
- **Fraudulent activities:** Intentional manipulation or misrepresentation.
- **Equipment malfunction:** Technical failures in data collection systems.

3 Why Outlier Detection Matters

3.1 Environmental Importance:

- **Protects Public Health:** Early warnings save lives during pollution spikes
- **Regulatory Compliance:** Keeps companies accountable to environmental laws

- **Urban Management:** Guide traffic controls and industrial activity during high pollution
- **Climate Research:** Identify extreme weather patterns and environmental shifts

3.2 Business & Economic Importance:

- **Financial:** Billions saved annually by catching fraud early
- **Quality Control:** Reduces waste and customer complaints
- **Enables Better Decisions:** Clean data leads to accurate predictions
- **Risk Management:** Identifies problems before they escalate

4 Key Detection Methods

4.1 IQR Method (Box Plot Rule)

- The Interquartile Range method provides a robust, non-parametric approach to outlier detection:

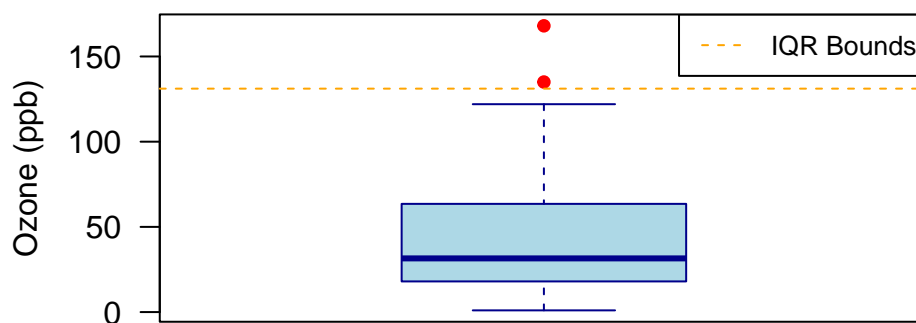
$$\text{Lower Bound} = Q1 - 1.5 \times IQR$$

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

$$\text{Outliers} = \{x_i \mid x_i < \text{Lower Bound} \vee x_i > \text{Upper Bound}\}$$

- Observations falling outside these bounds are classified as outliers. This method's robustness to extreme values makes it particularly suitable for environmental data analysis.

Ozone Distribution with Outliers



4.2 Z-Score Method

The Z-score standardizes each observation relative to the mean and standard deviation:

$$Z = \frac{X - \mu}{\sigma}$$

Values with $|Z| > 3$ typically indicate outliers, corresponding to observations beyond three standard deviations from the mean.

4.3 Local Outlier Factor (LOF)**

LOF compares the local density of an observation to that of its k nearest neighbors. Observations in sparse regions relative to their neighbors are flagged as outliers. This density-based approach excels at identifying contextual outliers that univariate methods miss.

5 Methodology: Implementing Outlier Detection in R

5.1 Dataset

The air quality dataset, available in the base R package, contains daily air quality measurements in New York from May to September 1973. The dataset includes 153 observations across 6 variables, with 111 complete cases after removing missing values Chambers et al. (1983).

5.2 Data Preparation and Exploration

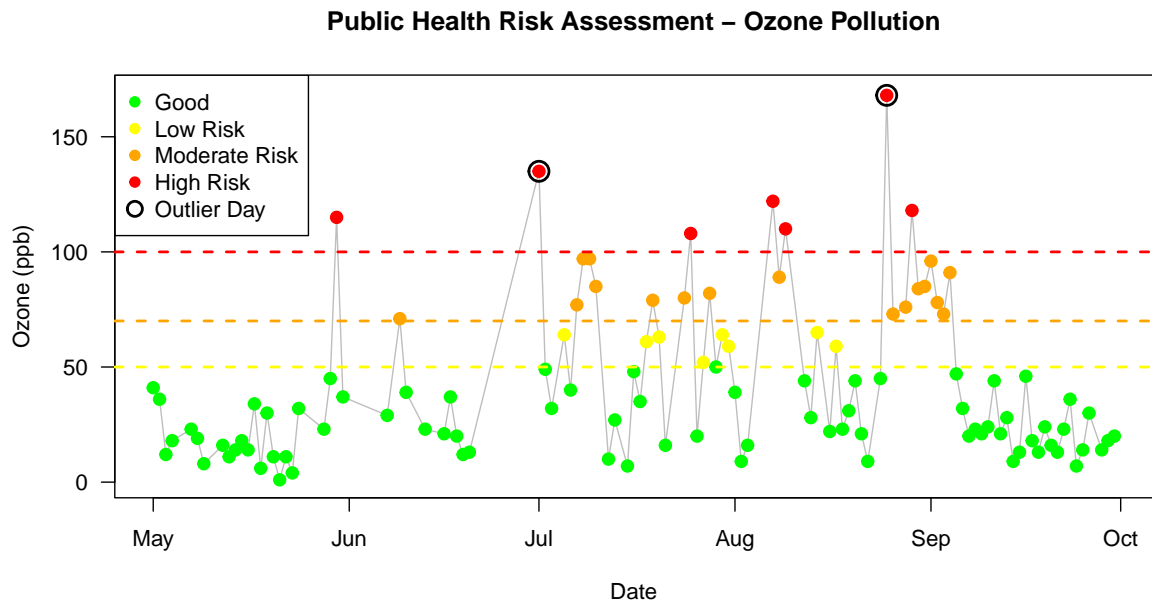
Ozone		Solar.R		Wind		Temp	
Min.	: 1.00	Min.	: 7.0	Min.	: 1.700	Min.	:56.00
1st Qu.:	18.00	1st Qu.:	115.8	1st Qu.:	7.400	1st Qu.:	72.00
Median	: 31.50	Median	:205.0	Median	: 9.700	Median	:79.00
Mean	: 42.13	Mean	:185.9	Mean	: 9.958	Mean	:77.88
3rd Qu.:	63.25	3rd Qu.:	258.8	3rd Qu.:	11.500	3rd Qu.:	85.00
Max.	:168.00	Max.	:334.0	Max.	:20.700	Max.	:97.00
NA's	:37	NA's	:7				

5.3 Ozone Outlier Analysis

Critical Pollution Event Detection — Ozone Outlier Analysis — IQR Bounds: -48 128 ppb
 Number of high ozone days: 2 Outlier ozone levels: 135 168 ppb

— Public Health Impact — Days with high health risk (Ozone > 100 ppb): 2 Days with moderate health risk (Ozone > 70 ppb): 0

5.4 Time Series Analysis of Pollution Events



6 Applications

6.1 Environmental Applications:

- **Public Health Monitoring:** Automated alerts when pollution exceeds safe levels.
- **Climate Research:** Identifying extreme weather and long-term climate shifts.

6.2 Business Applications:

- **Fraud Detection:** Spotting unusual financial transactions.
- **Sales Forecasting:** Using clean data to predict future revenue.

6.3 Technical & Research Applications:

- **Data Cleaning:** Removing errors before analysis.
- **Model Development:** Preparing data for machine learning.

7 Conclusion

This analysis demonstrated how IQR-based outlier detection in R can identify dangerous ozone levels in air quality data. Using the 1973 New York dataset, we found several days where ozone exceeded health safety thresholds, particularly during summer months when temperature and solar radiation drive ozone formation. The IQR method proved effective for environmental data as it doesn't require normal distribution and handles extreme values well. By combining statistical detection with public health thresholds, we transformed data analysis into actionable insights for environmental

protection. This approach provides a foundation for real-time monitoring systems that could help cities issue health advisories during high pollution events. As climate concerns grow, such methods will become increasingly valuable for protecting public health through data-driven environmental management.

8 Affidavit

I hereby affirm that this submitted paper was authored unaided and solely by me. Additionally, no other sources than those in the reference list were used. Parts of this paper, including tables and figures, that have been taken either verbatim or analogously from other works have in each case been properly cited with regard to their origin and authorship. This paper either in parts or in its entirety, be it in the same or similar form, has not been submitted to any other examination board and has not been published.

I acknowledge that the university may use plagiarism detection software to check my thesis. I agree to cooperate with any investigation of suspected plagiarism and to provide any additional information or evidence requested by the university.

8.1 Checklist

- ☒ The handout contains 7 pages of text.
- ☒ The submission contains the Quarto file of the handout.
- ☒ The submission contains the Quarto file of the presentation.
- ☒ The submission contains the HTML file of the handout.
- ☒ The submission contains the HTML file of the presentation.
- ☒ The submission contains the PDF file of the handout.
- ☒ The submission contains the PDF file of the presentation.
- ☒ The title page of the presentation and the handout contain personal details (name, email, matriculation number).
- ☒ The handout contains a bibliography, created using BibTeX with an APA citation style.
- ☒ Either the handout or the presentation contains R code that demonstrates coding expertise.
- ☒ The filled out Affidavit.
- ☒ The link to the presentation and the handout published on GitHub.

[Rajesh Mahato,] [10 Dec 2025,] [Bergheim]

9 References

- Chambers, J. M., Cleveland, W. S., Kleiner, B., & Tukey, P. A. (1983). *Graphical methods for data analysis*. Wadsworth & Brooks/Cole.
- Tukey, J. W. (1977). *Exploratory data analysis*. Addison-Wesley.