

# Comprehensive Anomaly Detection: Comparing Statistical and Machine Learning Methods

Multi-Method Outlier Identification Framework

---

## Problem Statement

Anomaly detection is critical across domains: fraud detection in finance, quality control in manufacturing, network intrusion detection, and sensor malfunction identification. However, no single method performs optimally across all scenarios. This project implements and compares multiple anomaly detection approaches—from classical statistical tests to modern machine learning algorithms—providing a comprehensive framework for identifying outliers under different data characteristics and business requirements.

## Business Value

Robust anomaly detection prevents financial losses from fraudulent transactions, identifies equipment failures before costly breakdowns, detects cybersecurity threats in real-time, and flags data quality issues in analytics pipelines. The multi-method approach ensures reliability across diverse operational contexts where false positives and false negatives have different cost implications.

## Methodology

**Dataset Characteristics:** Synthetic and real-world datasets with known outlier labels for method validation. Approximately 1,460 observations with injected anomalies representing 2-5% of data points, mimicking realistic anomaly rates in production systems.

### Method 1: Z-Score (Standard Deviation Method)

- Assumes normal distribution of data
- Flags points beyond 3 standard deviations from mean (99.7% confidence interval)
- Computationally efficient, interpretable threshold
- Effective for univariate outlier detection in normally distributed data
- Detected approximately 27 outliers (1.8% of dataset)

### Method 2: Interquartile Range (IQR Method)

- Distribution-free approach based on quartile spread
- Outlier threshold: values below  $Q1 - 1.5 \times IQR$  or above  $Q3 + 1.5 \times IQR$
- Robust to non-normal distributions and extreme values
- Widely used in exploratory data analysis and box plot visualization
- Identified approximately 37 outliers (2.5% of dataset)

### Method 3: Grubbs' Test (Extreme Deviation Test)

- Statistical hypothesis test for single outlier detection
- Compares maximum deviation from mean to critical value (t-distribution based)

- Iterative removal of most extreme outlier, retesting on remaining data
- Formal statistical significance testing (typical alpha = 0.05)
- Results confirmed absence of extreme single outliers in primary features

#### **Method 4: Isolation Forest (Ensemble Machine Learning)**

- Tree-based anomaly detection through random partitioning
- Anomalies isolated in fewer partitions (shorter path length in trees)
- Handles high-dimensional data and complex multivariate relationships
- Contamination parameter set to expected anomaly rate (0.05)
- Detected approximately 69 outliers (4.7% of dataset), capturing subtle patterns

#### **Method 5: DBSCAN (Density-Based Spatial Clustering)**

- Identifies outliers as points in low-density regions
- Parameters: epsilon (neighborhood radius), min\_samples (density threshold)
- Effective for non-linear cluster shapes and varying density regions
- Discovered approximately 22 outliers (1.5% of dataset) in isolated spatial regions

## **Results**

### **Comparative Analysis:**

- Detection rates varied from 1.5% (DBSCAN) to 4.7% (Isolation Forest)
- Isolation Forest captured most anomalies due to multivariate consideration
- Z-score and IQR showed moderate agreement (correlation 0.73)
- DBSCAN identified spatially isolated clusters missed by statistical methods
- Grubbs' test useful for initial screening but limited to univariate analysis

### **Method Selection Guidelines:**

- Z-score: Fast baseline for normally distributed features, real-time systems
- IQR: Robust choice for skewed distributions, financial data screening
- Isolation Forest: Best for high-dimensional data, complex multivariate anomalies
- DBSCAN: Spatial anomalies, geographic data, network topology analysis
- Ensemble approach: Combine multiple methods for high-stakes applications (e.g., fraud detection)

### **Production Deployment Insights:**

- Computational efficiency: Z-score (fastest) < IQR < DBSCAN < Isolation Forest (most comprehensive)
- For streaming data: Z-score with sliding window or incremental IQR
- For batch processing: Isolation Forest for maximum detection capability
- Recommend ensemble voting: flag point as anomaly if detected by 2+ methods

---

*This multi-method framework provides a comprehensive toolkit for anomaly detection, enabling practitioners to select optimal approaches based on data characteristics, computational constraints, and application requirements.*