# Efficient Data Stream Anomaly Detection

## Introduction

Anomalies, also known as outliers, are data points that deviate significantly from the norm. Detecting anomalies is crucial for identifying potential problems, unusual patterns, fraud, or significant events in various fields such as finance, cybersecurity, and healthcare. Anomalies are categorized into three main types:

1. **Point Anomalies**: Single data points deviating from the rest, like an unusually high credit card transaction.
2. **Contextual Anomalies**: Anomalies only within a certain context, such as a high temperature during winter.
3. **Collective Anomalies**: Groups of related data points showing abnormal behavior, such as a surge in network traffic.

**Note** - The green text code below shows a reference to the code

## Data Stream Simulation

The data stream simulation mimics real-world scenarios by generating a sequence of floating-point numbers, incorporating regular patterns, seasonality, and random noise. This approach helps test different anomaly detection techniques in a controlled setting.

**Code Explanation:**

The generate_data_stream function creates synthetic data with the following components:

- **Trend**: A slow upward drift (time * 0.0005), simulating a gradual increase over time.
- **Seasonal Pattern**: A sine function (2 * np.sin(2 * np.pi * time / seasonal_period)) represents repeating cycles, such as daily temperature changes or periodic financial trends.
- **Noise**: Random fluctuations (np.random.normal(0, 0.5, size=length)) to introduce natural randomness found in real data.
- **Anomalies**: Injected anomalies by randomly selecting data points and adding a large positive or negative value (e.g., ±10), simulating sudden spikes or dips.

The sine wave serves as a simple and interpretable representation of periodic data, making it ideal for validating different anomaly detection methods. Given this regular pattern, Z-Score detection performs well because it captures deviations from the mean effectively in a scenario where normal data behavior is straightforward.

# Anomaly Detection Algorithms

Four different algorithms were used: Z-Score, Holt-Winters, Isolation Forest, and Blending Ensemble with Hard Voting.
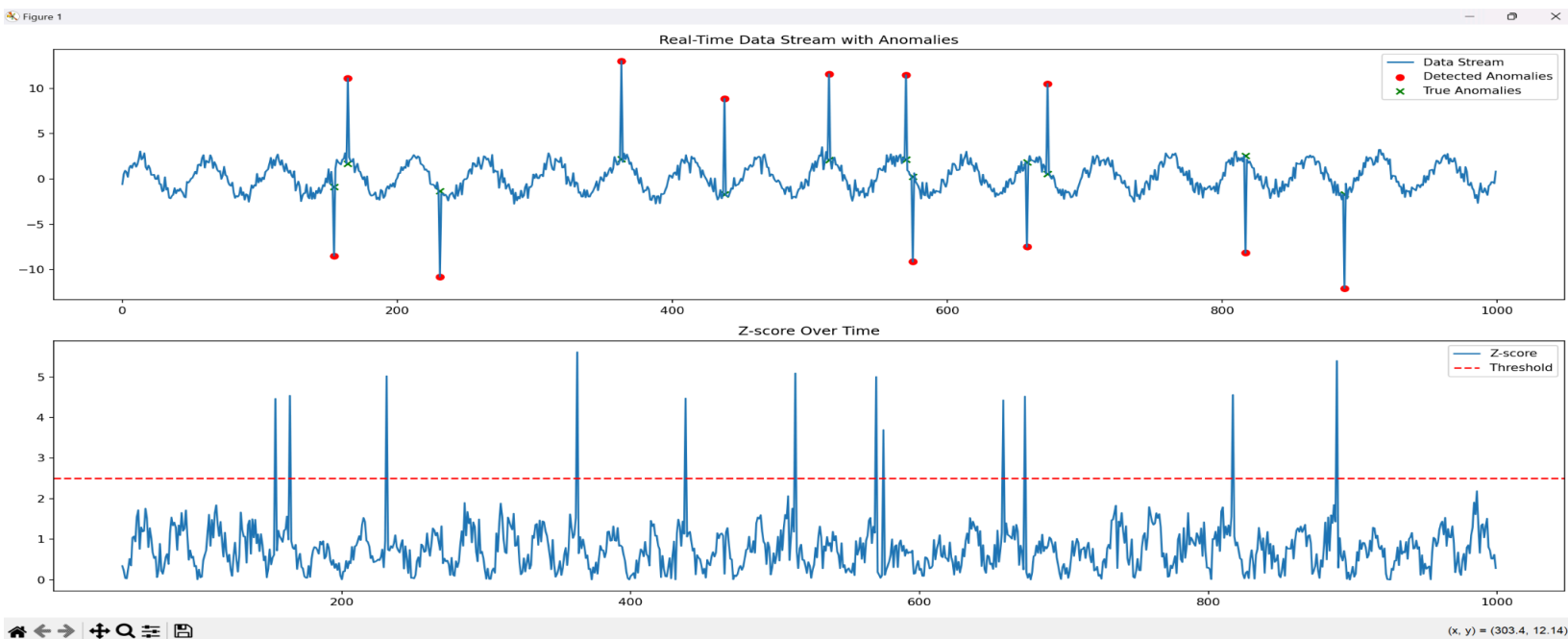
## 1. Z-Score Anomaly Detection

The Z-Score method measures how far each data point is from the mean of the windowed data in terms of standard deviations. It is suitable for simple patterns, such as the sine wave used here, due to the clear cyclical behavior, making it easy to detect deviations.

**Code Overview**:

- **Sliding Window Approach**: A fixed-size window (WINDOW_SIZE = 50) is used to calculate the rolling mean and standard deviation.
- **Thresholding**: Anomalies are identified if the absolute Z-score exceeds a specified threshold (THRESHOLD = 2.5).
- **Anomaly Detection**: For each new data point, the Z-score is calculated as (value - mean) / std, where mean and std are computed over the current window. If the standard deviation is zero (rare in real data but possible with constant values), it is adjusted to a small positive value to avoid division by zero.

The Z-Score technique performs well in this setup, achieving high accuracy because the sine function's regularity simplifies the detection of outliers.

# Graph(Z-Score Anomaly Detection) - From Code
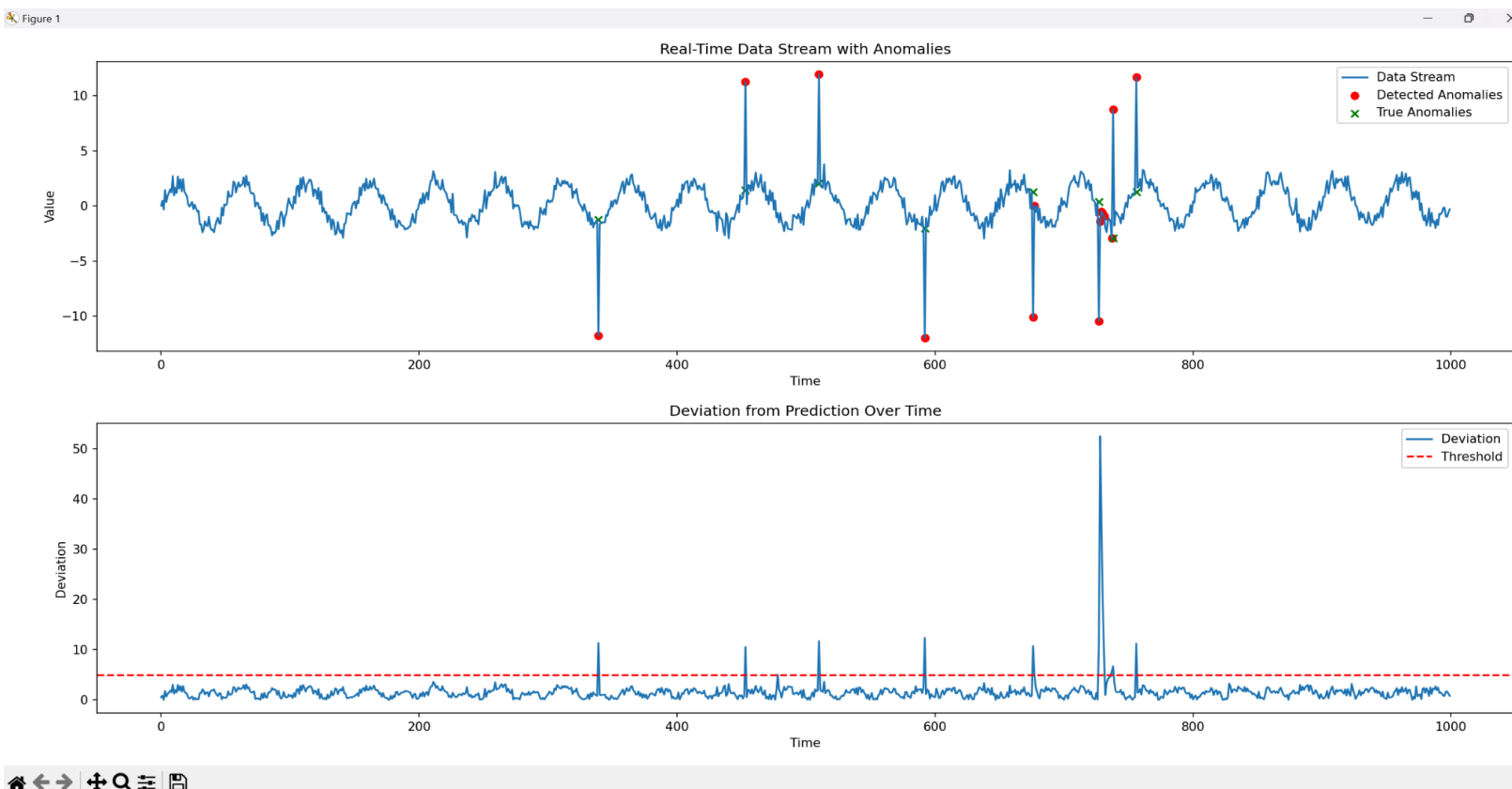
## 2. Holt-Winters Method

The Holt-Winters model applies exponential smoothing to account for level, trend, and seasonal components in time-series data, making it useful for more complex seasonal variations. It can predict the next value and detect anomalies based on deviations from this prediction.

**Code Explanation**:

- **Model Initialization**: The level (alpha), trend (beta), and seasonality (gamma) parameters control how much weight is given to new data points when updating the model.
- **Seasonal Decomposition**: The data is decomposed into trend, seasonal, and residual components, allowing the algorithm to adapt dynamically to changes in these patterns.
- **Prediction and Deviation Check**: For each incoming data point, the model predicts the expected value. If the difference between the actual and predicted value exceeds a threshold (THRESHOLD = 5.0), it is flagged as an anomaly.

The Holt-Winters method adapts well to changing patterns but may require fine-tuning of parameters to handle different seasonal behaviors effectively.

## Graph(Holt-Winters Method) - From Code
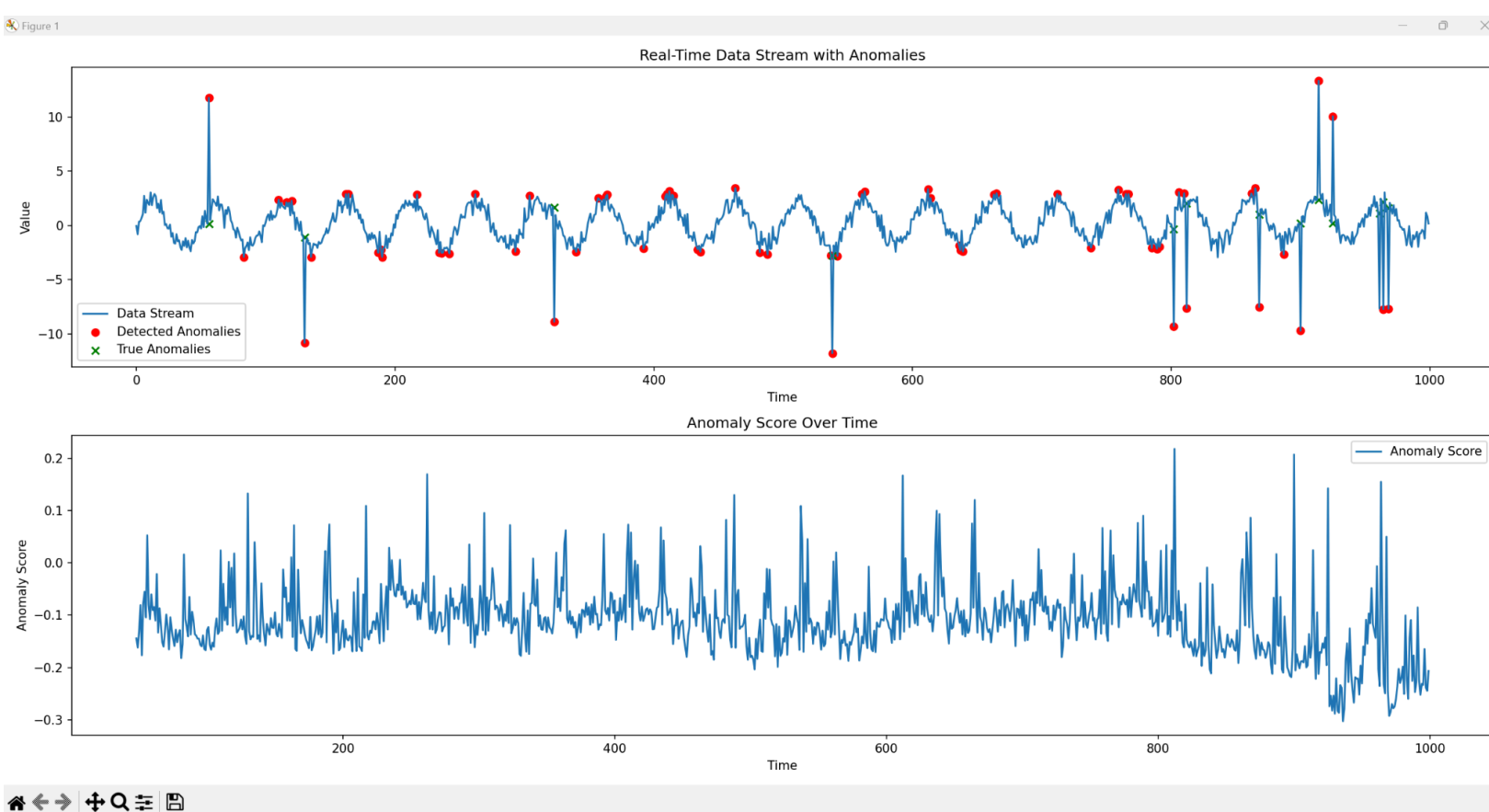
# 3. Isolation Forest

Isolation Forest identifies anomalies by isolating data points through random partitioning. It is well-suited for high-dimensional data and works effectively on continuous streams.

**Code Explanation**:

- **Sliding Window Training**: A window of data is used to update the Isolation Forest model periodically. The contamination parameter (CONTAMINATION = 0.05) estimates the proportion of anomalies in the data.
- **Anomaly Scoring**: The model assigns an anomaly score to each new data point. Scores above a certain level indicate potential anomalies.
- **Incremental Updates**: The model is retrained with each new window, allowing it to adapt to evolving data distributions (concept drift).

Isolation Forest handles varied data distributions but may struggle with highly regular patterns unless parameters are adjusted regularly.

# Graph(Isolation Forest) - From Code
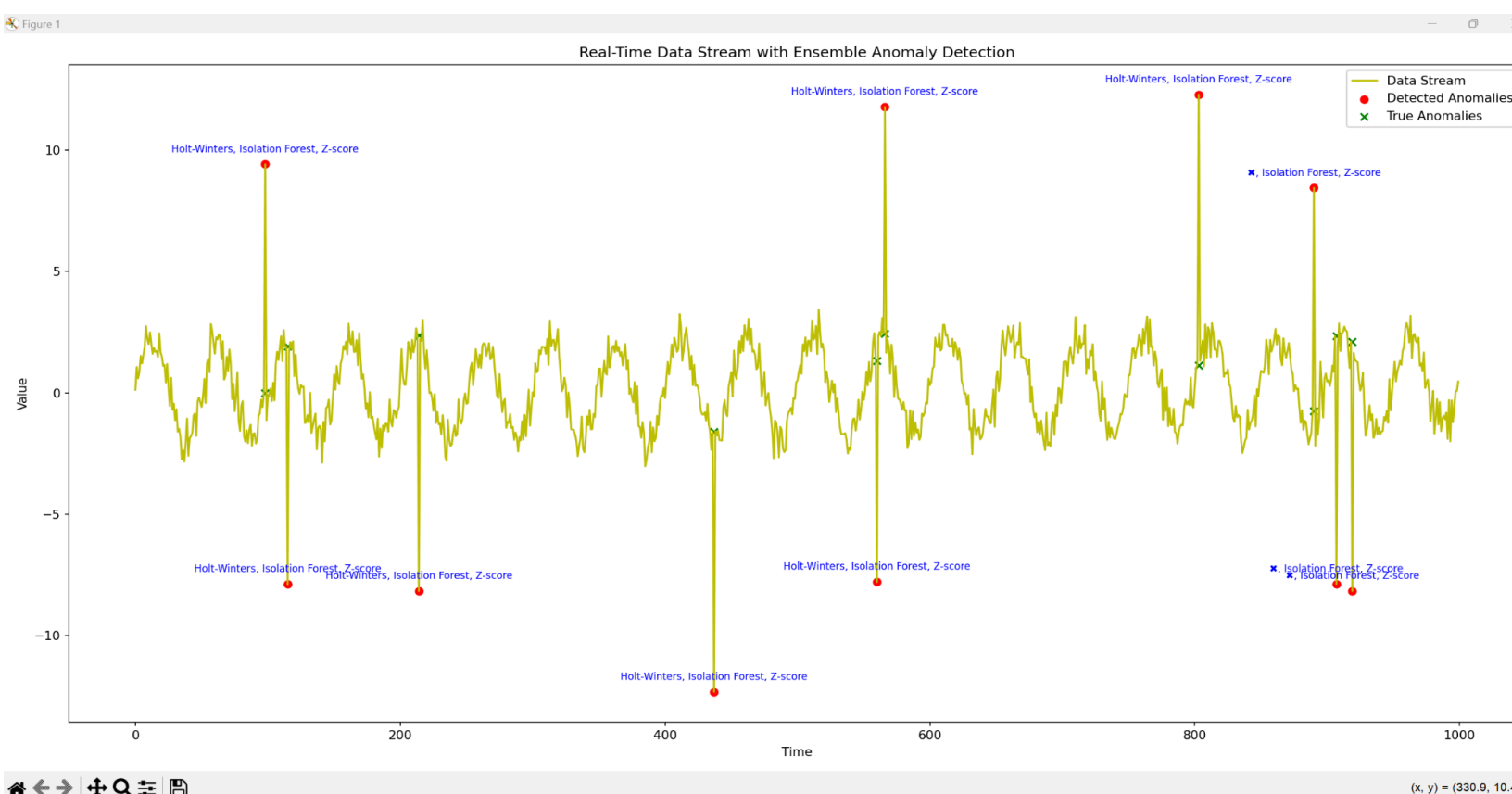
## 4. Blending Ensemble (Hard Voting)

The Blending Ensemble combines Z-Score, Holt-Winters, and Isolation Forest results using hard voting. It detects an anomaly when at least two of the three methods flag the same data point.

**Code Explanation**:

- **Ensemble Strategy**: For each new data point, predictions from Z-Score, Holt-Winters, and Isolation Forest are evaluated. The majority vote determines whether an anomaly is present.
- **Algorithm Annotation**: Annotations indicate which algorithms detected each anomaly, providing insights into detection consistency across methods.
- **Real-Time Adaptation**: The ensemble updates in real-time, incorporating input from all models and improving robustness to different types of anomalies.

The Blending Ensemble method improves detection accuracy by leveraging the strengths of each individual algorithm, achieving a balance between precision and recall.

# Graph(Blending Ensemble (Hard Voting)) - From Code

# Anomaly Detection Tracking and Annotation in Blending Ensemble

**Anomaly Detection Tracking:**

- It is checked whether the current data point is an anomaly detected by the ensemble.
- If an anomaly is found, the index and details of the algorithms (Holt-Winters, Isolation Forest, Z-score) that detected it are appended.
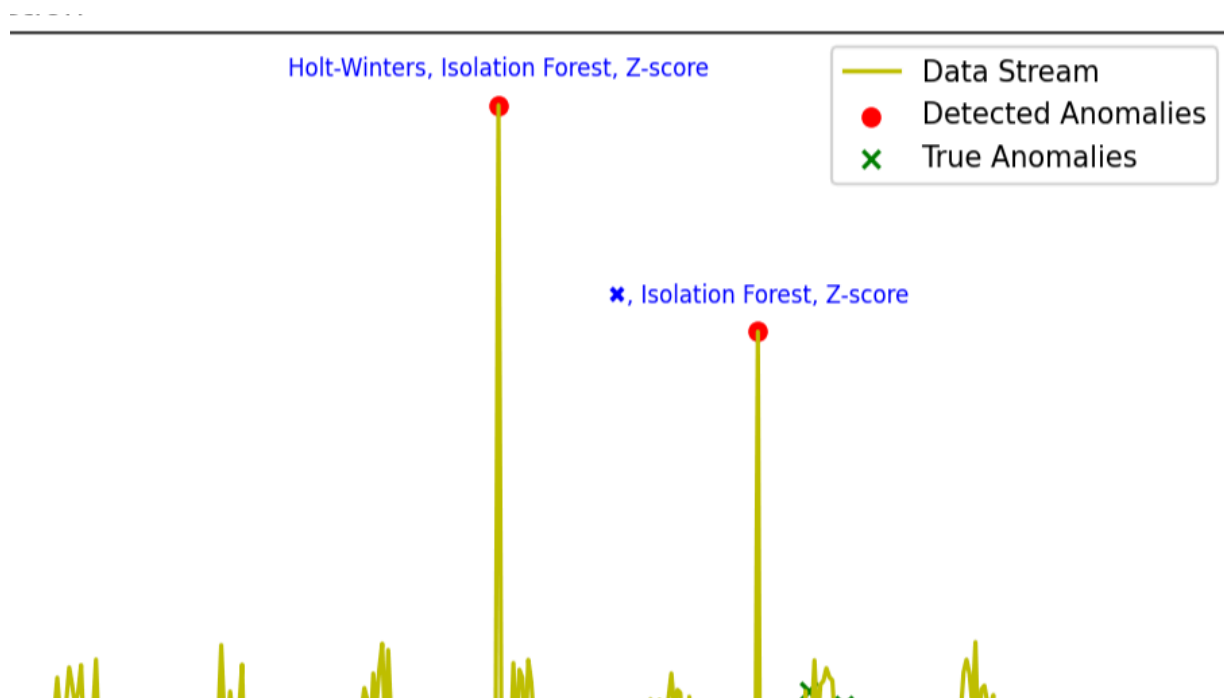
**Annotation Details:**

- The detected algorithms are marked with their names or '✖️' if they did not flag the anomaly.
- The details for each anomaly are stored.

**Adding Annotations:**

- For each detected anomaly, an annotation is added to the plot.
- The algorithms that flagged the anomaly are displayed in the annotation, positioned slightly above the data point in blue text.

**Example**

# Comparison

| Algorithm | Strengths | Limitations | Suitable Use Cases |
|-----------|-----------|-------------|--------------------|
| Z-Score | High precision for simple periodic data | Struggles with complex seasonal patterns | Regular data with clear cyclical behavior |
| Holt-Winters | Handles seasonal variations, adaptive to changes | Requires parameter tuning | Time series with known seasonal components |
| Isolation Forest | Effective for diverse data distributions | Needs periodic updates to manage concept drift | High-dimensional data, unstructured streams |
| Blending Ensemble | Combines multiple strengths, robust | Higher computational cost | Mixed data types with diverse anomaly patterns |

# Evaluation Matrix

| Metrics | Z Score | Holt-Winter | Isolation Forest | Ensemble |
|---------|---------|-------------|------------------|----------|
| Accuracy | 1.00 | 0.9500 | 0.9440 | 0.9990 |
| Precision | 1.00 | 0.1404 | 0.162 | 0.8889 |
| Recall | 1.00 | 0.8889 | 0.8462 | 1.000 |
| False Positive Rate | 0.00 | 0.0494 | 0.0547 | 0.0010 |
| False Negative Rate | 0.00 | 0.1111 | 0.1538 | 0.0000 |
| True Positive Rate | 1.00 | 0.889 | 0.8462 | 1.0000 |
| True Negative Rate | 1.00 | 0.9506 | 0.9453 | 0.9899 |
| F1 Score | 1.00 | 0.2424 | 0.2821 | 0.9412 |

# Conclusion

- **Z-Score** achieved the highest accuracy due to the regularity of the sine function used in data generation, which made deviations easier to detect.
- **Holt-Winters** effectively captured seasonal variations but was sensitive to parameter adjustments.
- **Isolation Forest** provided flexible detection across different data patterns but needed regular updates.
- **Blending Ensemble** delivered balanced performance by integrating multiple algorithms, making it more robust for varied anomaly detection needs.

# Future Work

Possible improvements include:

- Incorporating deep learning models for more complex, non-linear data patterns.
- Dynamic adjustment of algorithm parameters based on real-time feedback.
- Enhanced visualization capabilities for better monitoring and interpretation of anomaly detection results.