# Efficient Data Stream Anomaly Detection

## Introduction

This project aims to develop a robust Python script for detecting anomalies in a continuous data stream. Such streams could simulate various real-world metrics, including financial transactions, system performance metrics, or any other type of time-series data. Anomalies in these streams might include unusual spikes, drops, or deviations from the expected pattern. The main objectives of this project are to simulate a realistic data stream, implement an effective real-time anomaly detection algorithm, and provide visualizations to monitor the system's performance.

The project focuses on: - Simulating a data stream with normal and anomalous data points. - Detecting anomalies using an adaptive algorithm that adjusts to changes in the data. - Optimizing the detection algorithm for efficiency and accuracy. - Providing real-time visual feedback on the anomalies detected and overall data trends.

## Steps to Run the Code

1. **Set Up the Environment**: It is recommended to use a virtual environment to manage dependencies:
   - Create a virtual environment:
     ```
     python -m venv venv
     ```
   - Activate the virtual environment:
     - On Windows:
       ```
       venv\Scripts\activate
       ```
     - On macOS/Linux:
       ```
       source venv/bin/activate
       ```
2. **Install Dependencies**:
   - Install the required packages listed in `requirements.txt`:
     ```
     pip install -r requirements.txt
     ```
3. **Run the Script**:
   - Execute the main script to start the anomaly detection process:
     ```
     python main.py
     ```
4. **View Results**:
   - The script will generate real-time visualizations of the data stream and detected anomalies. Performance metrics and missed anomalies will be printed in the console for review.

## Algorithm Choice

### Adaptive Z-Score

**Reason for Choosing Adaptive Z-Score**: - The Adaptive Z-Score method is chosen due to its dynamic nature, which makes it well-suited for real-time

data streams where data distribution may change over time. This adaptability is crucial for accurately detecting anomalies in environments with evolving patterns or concept drift.

**Comparison with Other Algorithms**: - **Isolation Forest**: This algorithm isolates anomalies by randomly partitioning the data. While effective for high-dimensional data, it may not adapt quickly to changes in data patterns over time, which is crucial for real-time detection. - **One-Class SVM**: This approach creates a decision boundary around normal data and identifies outliers outside this boundary. However, it is computationally intensive and less practical for large, continuous data streams due to its high complexity and resource requirements. - **LSTM-based Approaches**: Long Short-Term Memory networks are powerful for handling sequence data and learning temporal dependencies. However, they require substantial computational resources and training data, making them less suitable for real-time applications compared to the Adaptive Z-Score.

**Advantages of Adaptive Z-Score**: - **Balance of Accuracy and Efficiency**: The Adaptive Z-Score method offers a good balance between detection accuracy and computational efficiency. It adjusts its parameters based on recent data, allowing it to detect anomalies promptly while maintaining reasonable processing times. - **Adaptability**: It effectively handles changes in data trends and seasonal variations by continuously updating its mean and standard deviation estimates.

## Limitations

### Warm-Up Period

- The Adaptive Z-Score method requires a warm-up period to establish initial threshold values, which may affect the detection of early anomalies. During this phase, the threshold gradually adjusts, potentially leading to reduced sensitivity to anomalies until stabilization is achieved.

### Other Limitations

- **Parameter Sensitivity**: The effectiveness of the Adaptive Z-Score method can be influenced by the choice of parameters, such as the smoothing factor (alpha) and Z-score thresholds. Incorrect parameter settings may lead to either too many false positives or missed anomalies.
- **Handling Highly Volatile Data**: The method may struggle with extremely volatile data without proper tuning, as excessive noise can impact the calculation of mean and standard deviation.

## Mitigations and Fixes

### Warm-Up Period Adjustment

- **Tuning the Warm-Up Period**: The warm-up period has been optimized through experimentation to minimize its impact on early anomaly detection. The current settings provide a balance between initialization and sensitivity, reducing the risk of missing early anomalies.

### Parameter Tuning

- **Parameter Optimization**: Extensive testing was conducted to find optimal values for parameters like alpha, initial Z-score threshold, and final Z-score threshold. This fine-tuning helps in balancing detection sensitivity and minimizing false positives.

**Pros of These Fixes**: - **Improved Early Detection**: By optimizing the warm-up period and parameter values, the system can more accurately detect anomalies from the start, enhancing overall performance. - **Enhanced Accuracy**: The adjustments lead to better accuracy in identifying true anomalies while reducing the rate of false positives.

## Conclusion

The "Efficient Data Stream Anomaly Detection" project demonstrates an effective approach to detecting anomalies in continuous data streams using the Adaptive Z-Score method. The project successfully integrates a real-time data stream simulation with a robust anomaly detection algorithm and provides visual tools for monitoring and analysis. While there are inherent limitations, particularly related to the warm-up period and parameter sensitivity, these have been addressed through careful tuning and optimization. The result is a system that offers a strong balance between accuracy and efficiency, making it suitable for various real-time applications.