

# Open-source Technologies and Stream Mining joint Project Documentation

## Table of contents

<b>1</b>	<b>Introduction</b>	<b>3</b>
<b>2</b>	<b>Background and Literature Review</b>	<b>3</b>
<b>3</b>	<b>Dataset</b>	<b>3</b>
<b>4</b>	<b>System architecture</b>	<b>4</b>
<b>5</b>	<b>Modeling and Predictions</b>	<b>5</b>
5.1	Overview . . . . .	5
5.2	Anomaly detection . . . . .	5
<b>6</b>	<b>Experiments and testing</b>	<b>5</b>
<b>7</b>	<b>References</b>	<b>6</b>

*Smart City Air Quality Monitoring with Real-Time Stream Analytics (SCAir-IoT)*

**Abstract**

lorem ipsum...

**Keywords:** *Anomaly detection, Forecasting, Stream mining, Open-source technologies*

**Authors**

Albazzal, Houmam  
Balogh, Máté  
Földvári, Ádám  
Lahmar, Abderraouf  
Nagy, Zsuzsanna

**Mentored by**

Loubna Seddiki

# 1 Introduction

## 2 Background and Literature Review

In stream mining, we are limited to a portion of data and make decisions real-time in memory. As *Wares, Isaacs, and Elyan (2019)* highlight, in traditional machine learning contexts, data is referred to as batch data which can be loaded into memory in its entirety. According to the authors,

“this is of stark contrast to stream mining, where data streams produce elements in a sequential, continuous fashion, and may also be impermanent, or transient, in nature ... This means stream data may only be available for a short time.”

The authors refer to *Babcock et al. (2002)*, highlighting that

“once an element from a data stream has been processed, it is discarded or archived. It cannot be retrieved easily unless it is explicitly stored in memory, which is small relative to the size of data streams.”

## 3 Dataset

The dataset is the *UCI Air Quality* dataset *Vito, S. (2008)* which includes responses of gas sensor devices deployed in an Italian city. Besides these device readings, each gas measurement has a counterpart feature which denotes the gas concentration recorded by a co-located certified analyzer. Additionally, readings related to temperature along with absolute and relative humidity are included in the dataset.

The records span 1 year from March 2004 to February 2005, and present hourly aggregated measurements. Missing values are denoted with the value of -200.

## 4 System architecture

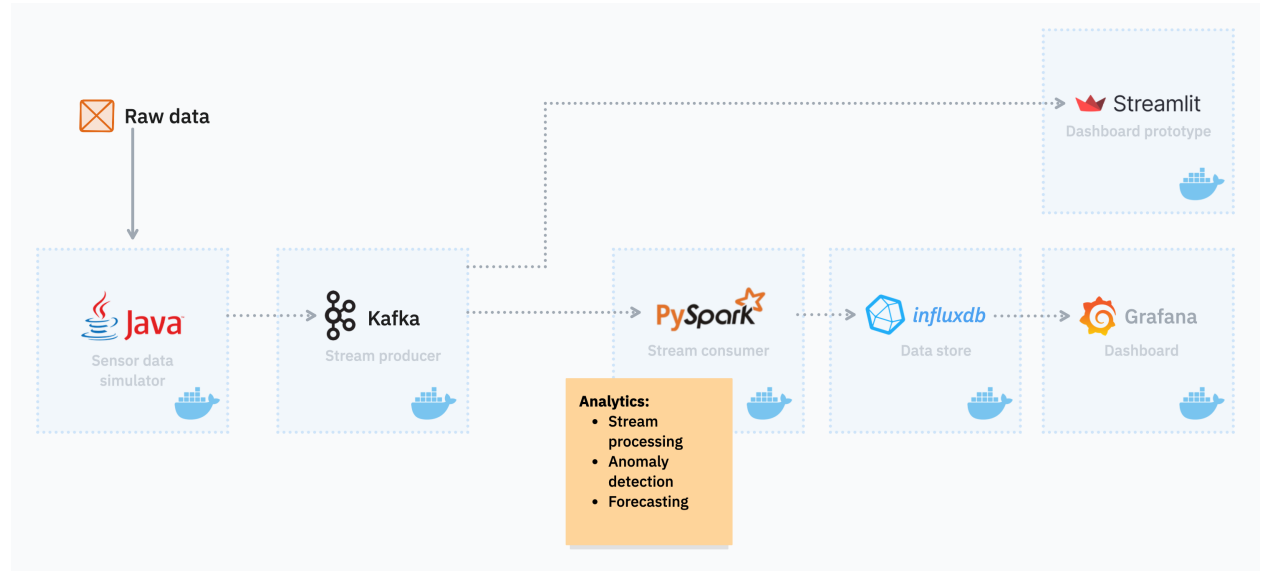


Figure 1: High-level view of the architecture with the utilized open-source technologies denoted for each component.

The stream mining pipeline includes components of *simulator*, *producer*, *consumer*, *data store*, and *dashboards* by which the raw `csv` dataset file was ingested. We utilized open-source technologies of *Java*, *Kafka*, *PySpark*, *Influxdb*, *Streamlit*, and *Grafana* provisioned in a containerized environment via *Docker*. The responsibility of each component is summarized as follows:

**Raw data:** Static file containing sensor measurements related to air quality.

**Sensor data simulator:** Reads the raw data file and simulates flow of sensor data.

**Kafka producer:** Streamlines the simulated sensor data into Kafka topics in their datetime order.

**PySpark consumer:** Listens to the streamlined data and creates mini batches to call analytical functions—such as Anomaly Detection and Forecasting—on this windowed data.

**Influxdb:** Data is then persisted in the database including the online predictions and the original incoming data.

**Streamlit:** Streamlit is used to create dashboard prototypes without the necessity of a database storage and connection. Kafka messages are directly consumed by this component to display simple line charts and anomaly detection related information and alerts.

**Grafana:** Dashboard visualization component that periodically fetches the database for new data to show the latest insights in real-time.

## 5 Modeling and Predictions

### 5.1 Overview

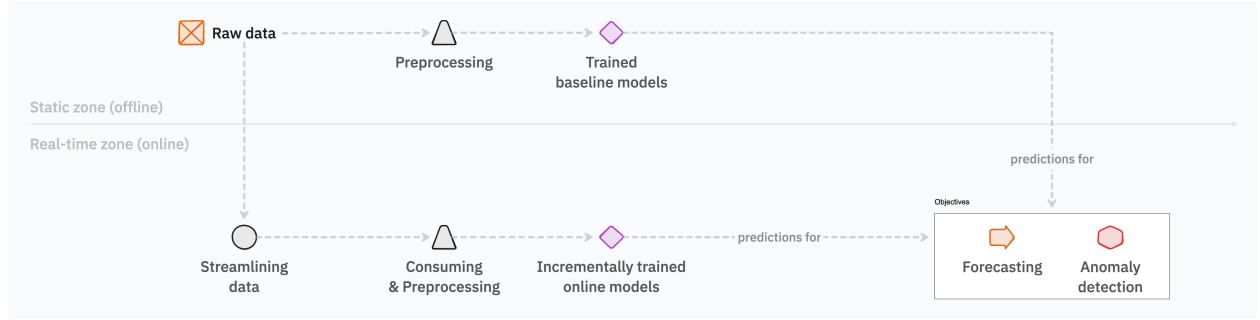


Figure 2: High-level view of the offline and online machine learning pipeline.

### 5.2 Anomaly detection

For *Anomaly Detection*, different techniques were utilized ranging from recognizing sensor readings with missing measurement to detecting local and global outliers, in each sensor separately.

Both, local and global approaches were developed with the general novelty, outlier detection methodology in mind, where the steps are:

1. Transform the data, to better reflect properties of interest
2. Obtain a novelty function
3. Apply peak picking algorithm on the novelty function

To detect local patterns in form of sudden changes (peaks and valleys), an 8-sample window was utilized on the streamlined data. As part of the data transformation, values of  $-200$  were replaced with **NA**. The most recent element in the window was dedicated as a test sample to make predictions for, while the rest of the in-window samples (at most 7) were designated as historical training points. Derivative operation was applied on the window data to turn a measurement value  $x_T$  into the change measured from  $x_{T-1}$ . This resulted in the novelty function, where **NA** values were first imputed with the window median but were transformed back to **NA** once the difference values were obtained. During the “training”, the in-window estimator calculates statistics ( $IQR, Q1, Q3$ ) from the in-window historical samples to perform a traditional outlier detection test using the  $Q1 - 1.5 \times IQR$  and  $Q3 + 1.5 \times IQR$  as thresholds to flag the test sample as outlier.

To account for out-of-window global patterns and to implement a detector that is applicable in stream mining systems where we cannot always rely on storing every incoming data point, an online detector was implemented using the *T-digest* algorithm. This method uses the same peak picking strategy using the traditional outlier thresholds, yet the quartiles along with the *Interquartile Range* get iteratively updated by each new sample. The algorithm maintains a centroid-based representation about these statistics using the incoming data and updates its state through merges.

## 6 Experiments and testing

TODO: add outputs of anomaly detection

## 7 References

Wares, S., Isaacs, J. and Elyan, E. (2019). Data stream mining: methods and challenges for handling concept drift. *SN Applied Sciences*, 1(11), 1412.

Babcock, B., Babu, S., Datar, M., Motwani, R. and Widom, J. (2002, June). Models and issues in data stream systems. In *Proceedings of the twenty-first ACM SIGMOD-SIGACT-SIGART symposium on Principles of database systems* (pp. 1–16).

Vito, S. (2008) Air Quality Dataset. UCI Machine Learning Repository. Available at: <https://archive.ics.uci.edu/ml/datasets/Air+Quality> (Accessed: 25 November 2025).

TODO: reference *T-digest*