

Real-Time Anomaly Detection with Apache Beam

Shunping Huang



Anomaly Detection



- Anomaly: the data instance that is **different** from the **normal** ones.
- Anomaly Detection: the task to identify the anomalies.

Anomalies in Different Types of Data

• Bounded vs. Unbounded





• Univariate vs. Multivariate





Anomalies in Different Types of Data (cont.)

• Data Streams vs. Time Series





• Data Streams with Concept Drift



One-Pager: Beam Basics

- A framework to unify the **batch** (bounded) and **streaming** (unbounded) processing.
- Key Concepts
 - PCollection a representation of data for parallel processing
 - **PTransform** a representation of computation to transform data
 - **Pipeline** a Directed Acyclic Graph (DAG) of PTransforms



PTransform for Anomaly Detection



Internals

• Individual Detectors

- Run a specific anomaly detection method and generate scores or labels
- e.g. <u>Incremental Z-Score</u>, <u>IQR</u>, and offline models (such as <u>isolation forests</u>, <u>LOF</u>, <u>one-class SVM</u>, etc) supported by PyOD.

• Threshold Criteria

• Fixed Thresholding, Quantile Thresholding



Internals (cont.)

• Ensemble Detectors

• Run a set of sub-detectors in parallel and aggregate scores or labels

• Aggregation Strategies

- $\circ~$ For Scores: Average Score, Max Score
- For Labels: Majority Vote, All Vote, Any Vote



Demo

Key Advantages of Beam for Anomaly Detection





• Design Doc

https://docs.google.com/document/d/1tE8lz9U_vjlNn2H7t-GRrs3vfhQ5UuCgWiHXCRHRPns/edit?usp =sharing

Source Code

https://github.com/apache/beam/tree/master/sdks/python/apache_beam/ml/anomaly

• Python Doc

https://beam.apache.org/releases/pydoc/current/apache_beam.ml.anomaly.html

• Colabs

https://github.com/apache/beam/tree/master/examples/notebooks/beam-ml/anomaly_detection

Thank you!

Try our new Apache Beam anomaly detection and give us feedback. Contributions are very welcome!

Shunping Huang shunping@google.com https://www.linkedin.com/in/shunping-huang-545b0150

