# Time Series Anomaly Detection: An Overview

Paul Boniol

Inria, ENS, PSL University

Paul.boniol@inria.fr

# Introduction: *Time series are Everywhere*

| Energy Production | Astrophysics | Medicine | Volcanology |
|---|---|---|---|



Edf.fr: tinyurl.com/yc7x5xje

Virgo: https://www.virgo-gw.eu/

tinyurl.com/39dx2us4

tinyurl.com/ybcttmfz

# Introduction: *Time series are Everywhere*

| Energy Production | Astrophysics | Medicine | Volcanology |
|---|---|---|---|

Secondary circuit sensor measurements
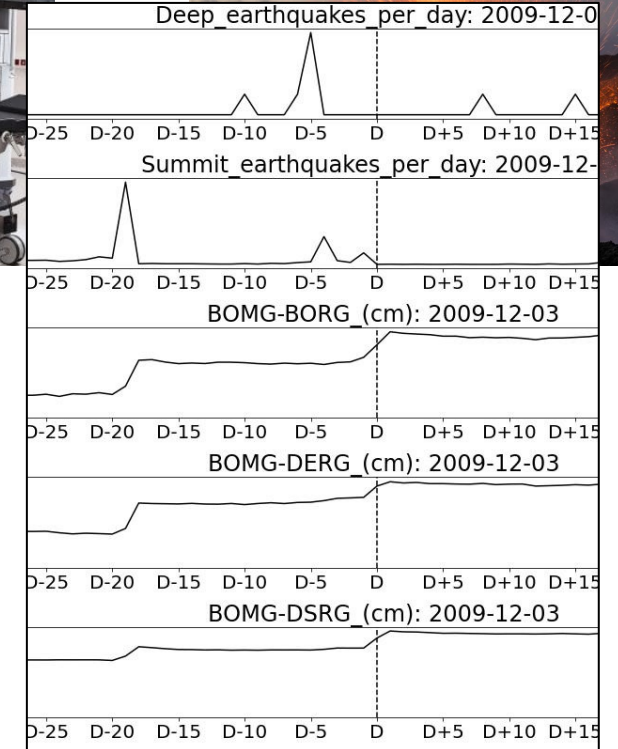
Fiber-acoustic sensors in the VIRGO north building

Sensor measurements of the Da-Vinci surgery robot

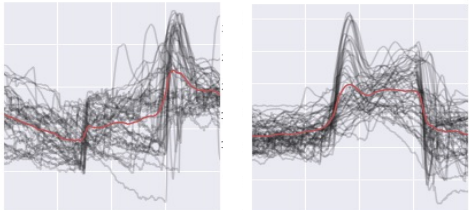Sensor measurements on le Piton de la Fournaise

# Introduction: *with Important Challenges*

**Energy Production**

Secondary circuit sensor measurements

Identification of precursors of feed-water pumps vibrations

**Astrophysics**

Fiber-acoustic sensors in the VIRGO north building

Noise detection in VIRGO interferometer north building

**Medicine**

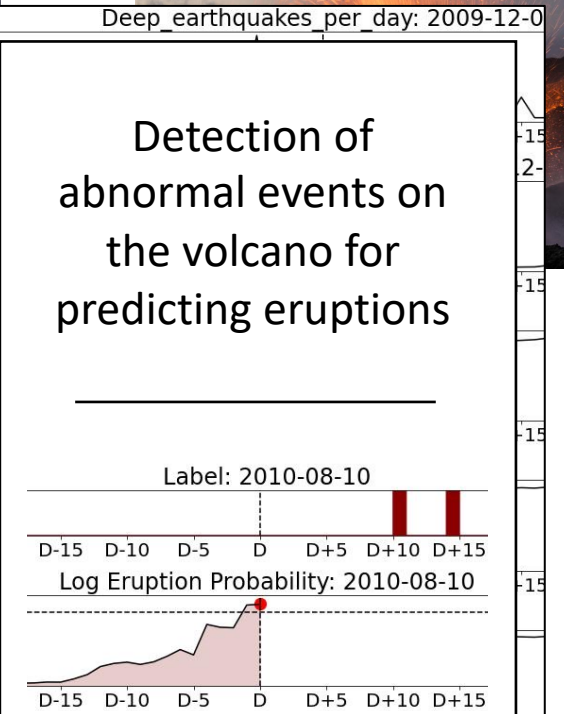Sensor measurements of the Da-Vinci surgery robot

Unusual surgeons gestures detection

**Volcanology**

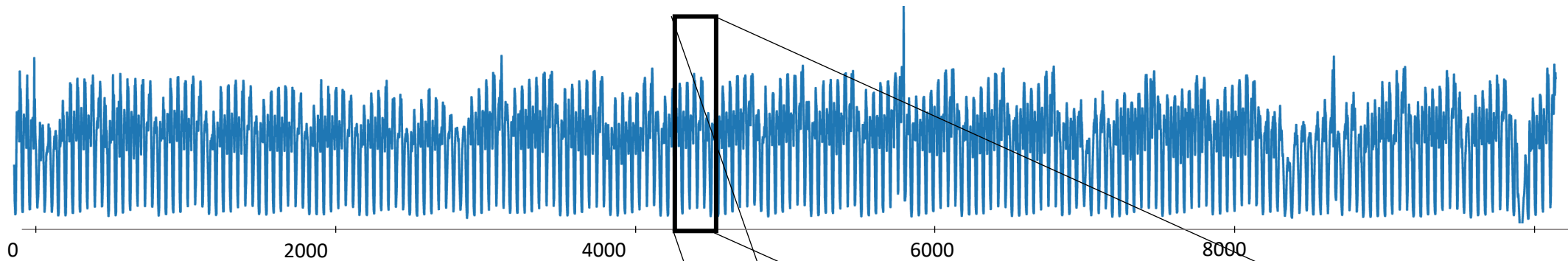Sensor measurements on le Piton de la Fournaise

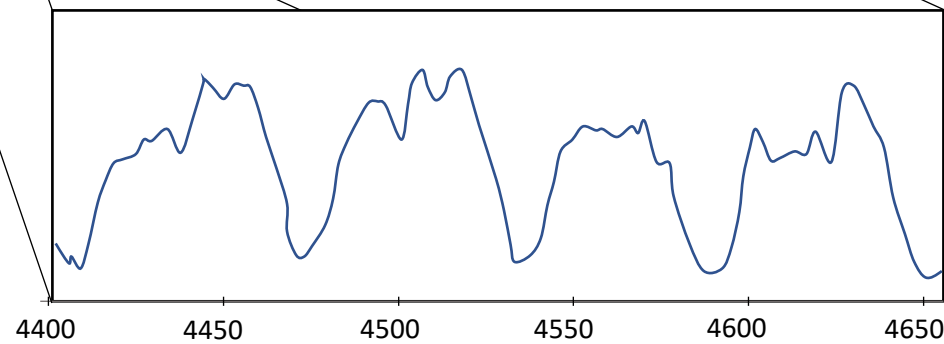Detection of abnormal events on the volcano for predicting eruptions

# Introduction: *Anomaly Detection in Time Series*

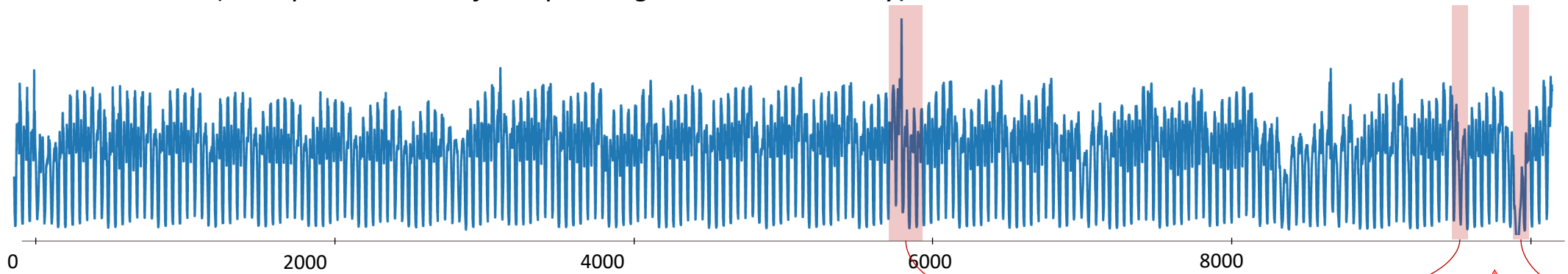- Time series $T$ (example : number of taxi passengers in New York City)



- Subsequence $T_{i,\ell}$
  with $i = 4400, \ell = 250$

# Introduction: *Anomaly Detection in Time Series*

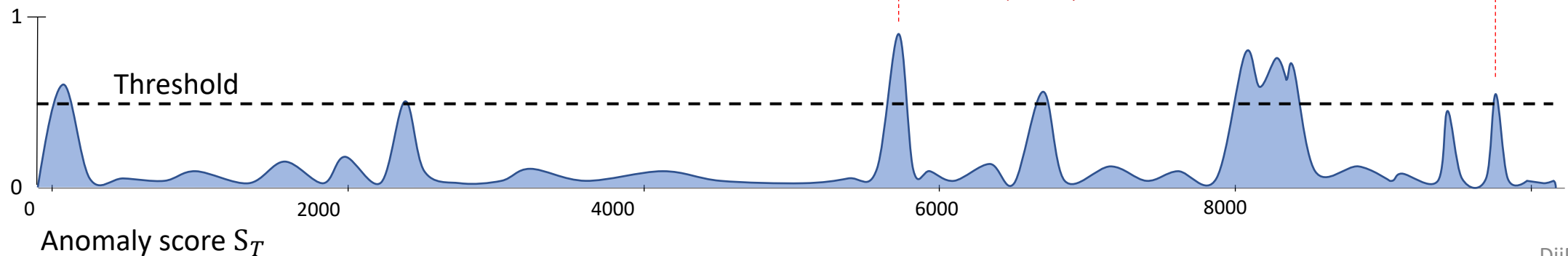- Time series $T$ (example : number of taxi passengers in New York City)

- *Anomaly: <span style="color:red">rare</span> point or sequence (of a given length) potentially <span style="color:red">non-desired</span>*
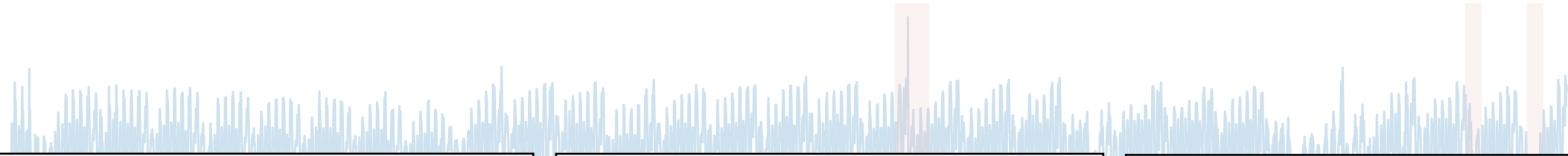
Daylight Saving Time (DST)

Flooding

Snowstorm

Threshold

Anomaly score $S_T$

# Introduction: *Outline*

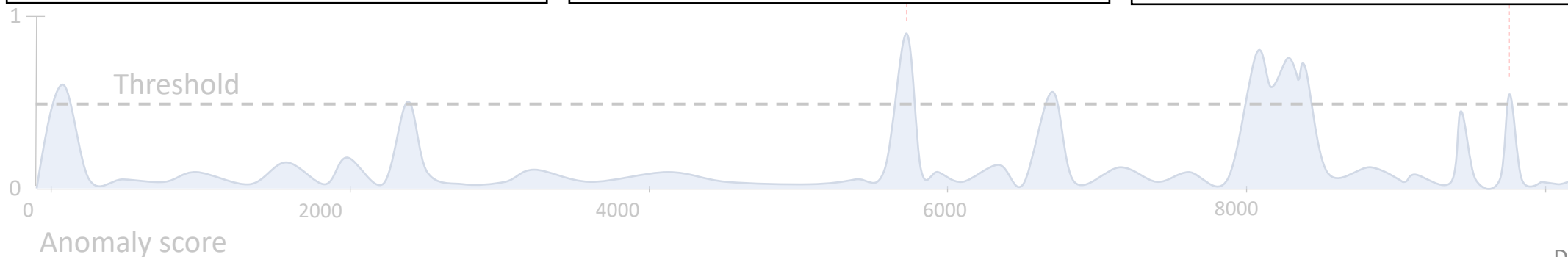- Time series *(example : number of taxi passengers in New York City)*



## 1. Foundations

1.1. Type of Time Series
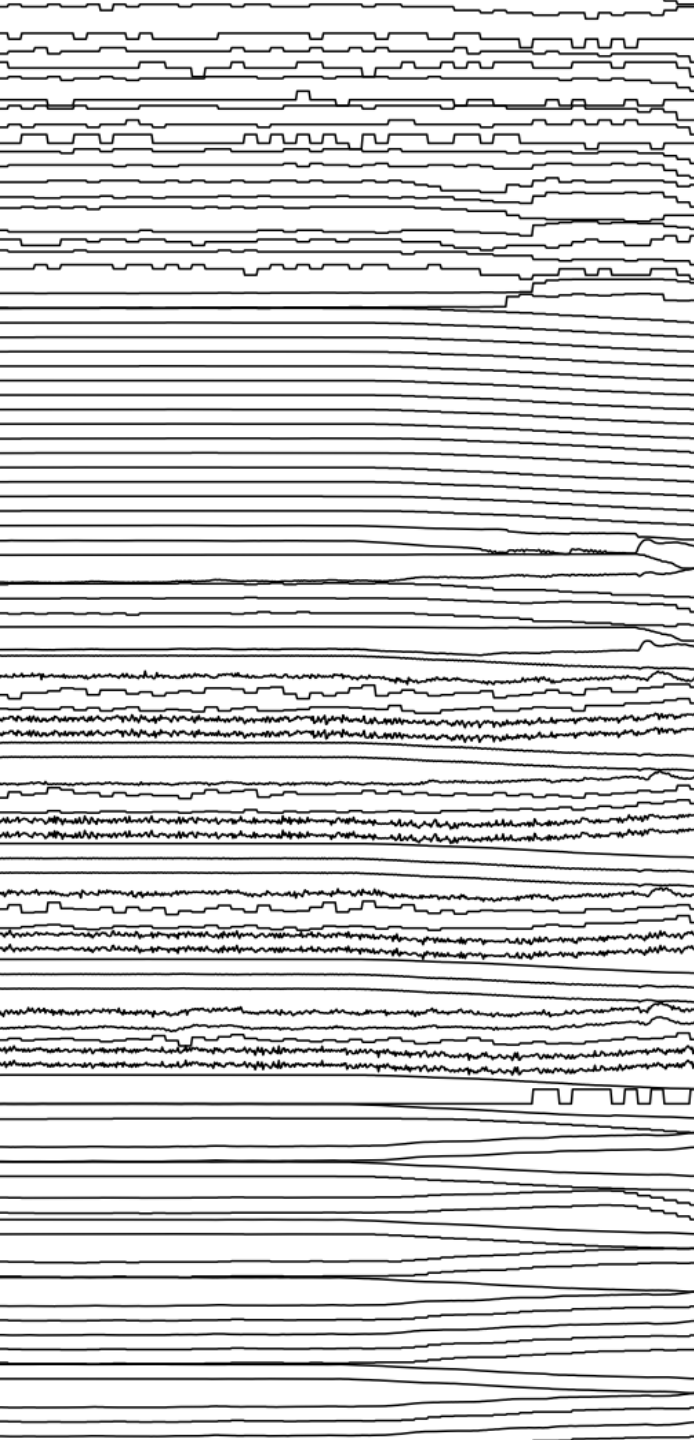1.2. Type of Anomalies

## 2. Anomaly Detection Methods

2.1. A Taxonomy of Methods
2.2. Existing Benchmarks

## 3. Perspectives and challenges

3.1. Time series labeling issues
3.2. Ensembling

# Foundations
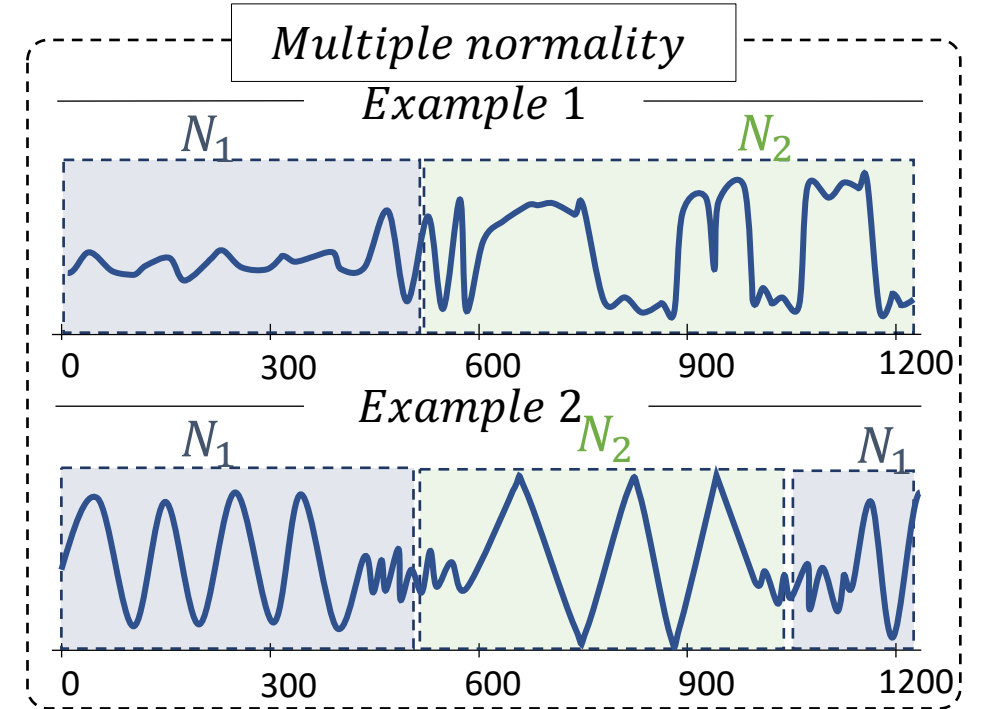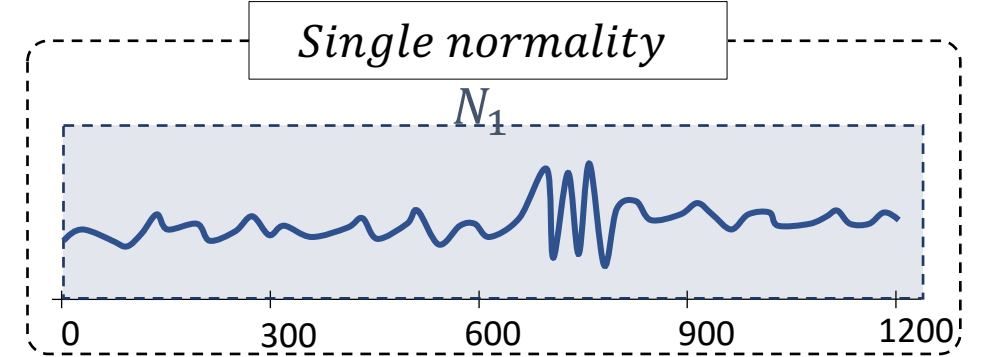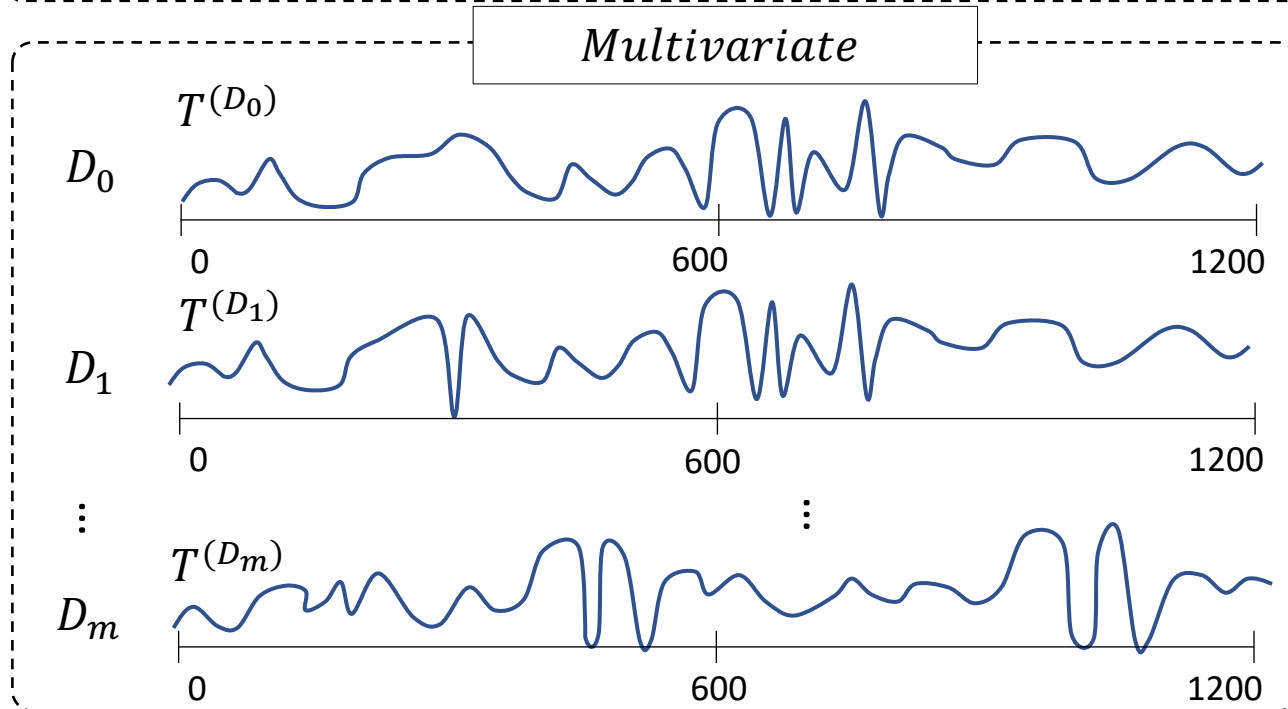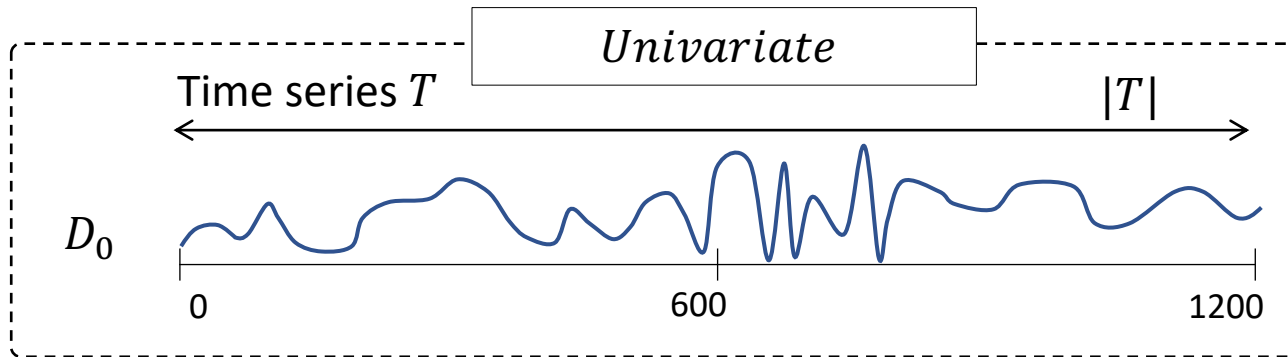
# Foundations: *Type of time series*

# Foundations: *Type of anomalies*

Example of single anomaly [3]

(a.1) Point-based

Time series — point anomaly

Contextual anomaly

distribution

distribution    distribution

(a.1.1) Point anomaly outside the healthy range of values (dotted black line)

(a.1.2) Contextual anomaly outside the local healthy range of values (dotted black box)

(a.2) Subsequence-based

Collective anomaly    Time series

distribution

(a.2) collective anomaly inside the healthy range of values (dotted black line)

(b.1) Single anomaly

$A_1$

(b.2) Multiple anomaly

(b.2.1) Multiple *different* anomaly

$A_2$    $A_1$
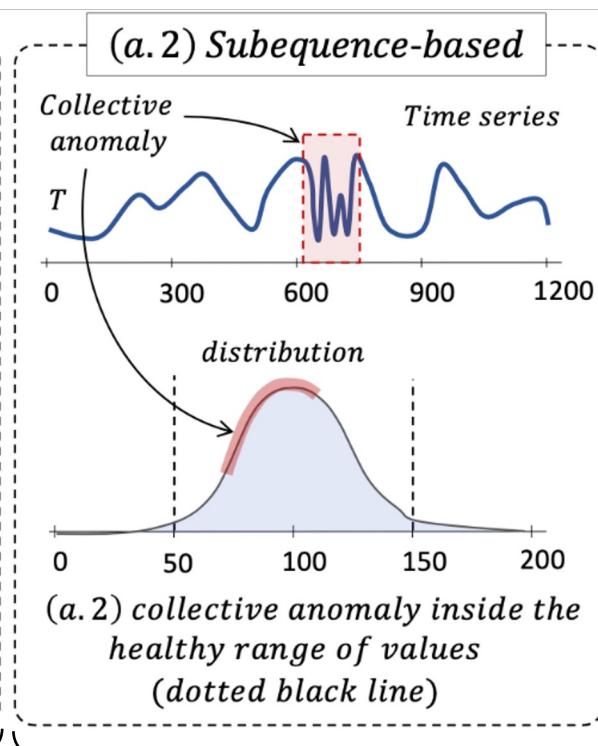
(b.2.2) Multiple *similar* anomaly

$A_1$    $A_1$

Example of point-based anomaly [1]

Example of subsequence-based anomaly [2]

Example of multiple anomaly [4]

# Foundations: *Type of anomalies*



Univariate and **Multivariate** point anomalies

Univariate
value outlier
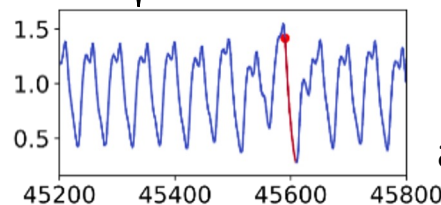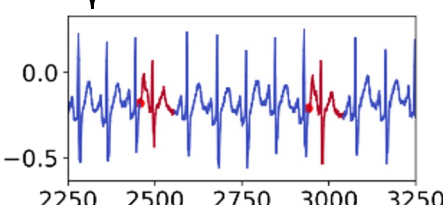
Multivariate
value outlier

(a.1) *Univariate case*

(a.2) *Multivariate case*

(a) *Point outlier outside the healthy range of values*
*(dotted black line)*

Univariate and **Multivariate** sequence anomalies

Univariate
subsequence outlier

Multivariate
subsequence outlier

(b.1) *Univariate case*

(b.2) *Multivariate case*

(b) *Subsequence outlier inside the healthy range of values*
*(dotted black line)*

# Anomaly Detection Methods

# Anomaly Detection methods: *A taxonomy*



Anomalies

Time Series

Pre-Processing

Anomaly Detection Model

Scoring

Post-Processing

Anomaly score

# Anomaly Detection methods: *A taxonomy*

By domains [5] …



[5] Sebastian Schmidl, Phillip Wenig, and Thorsten Papenbrock. 2022. Anomaly detection in time series: a comprehensive evaluation. Proc. VLDB Endow. 15, 9 (May 2022), 1779–1797.

# Anomaly Detection methods: *A taxonomy*

By inputs…



Time series anomaly detection methods

**Supervised**

**Semi-supervised**

**Unsupervised**

Training dataset
- Normal examples
- Anomaly examples

Time Series T

Training dataset
- Normal examples

Time Series T

Time Series T

# Anomaly Dete

## By inputs...

Time

**Supervised**

Training dataset
- *Normal examples*
- *Anomaly examples*

*Time Series T*



VVP: main steam system
28 sensors
*(flow, pressure, temperature)*

GSS: moisture separator-reheater system
2 sensors *(pressure and temperature)*

CEX: condensate extraction system
2 sensors *(pressure and temperature)*

GCT: turbine bypass system
2 sensors *(pressure)*

KKO: energy metering system
1 sensor *(power)*

*Primary circuit*

*Secondary circuit*

CEX

VVP — *Steam*

GRE

GCT

Turbine

KKO

GSS

ASG

ADG

*High pressure water*

Steam generator

ARE ← AHP

*Feed-water Pumps*

Condensor

*Cold water*

# Anomaly Dete...

By inputs...

Supervised

VVP: main steam system
28 sensors
*(flow, pressure, temperature)*

CEX: condensate extraction system
*2 sensors (pressure and temperature)*

GSS: moisture separator-reheater system
*2 sensors (pressure and temperature)*

GCT: turbine bypass system
*2 sensors (pressure)*

KKO: energy metering system
*1 sensor (power)*

*Primary circuit*

*Secondary circuit*

CEX

VVP

*Steam*

GSS

GCT

GRE

Turbine

KKO

ASG

ADG

Steam generator

*High pressure water*
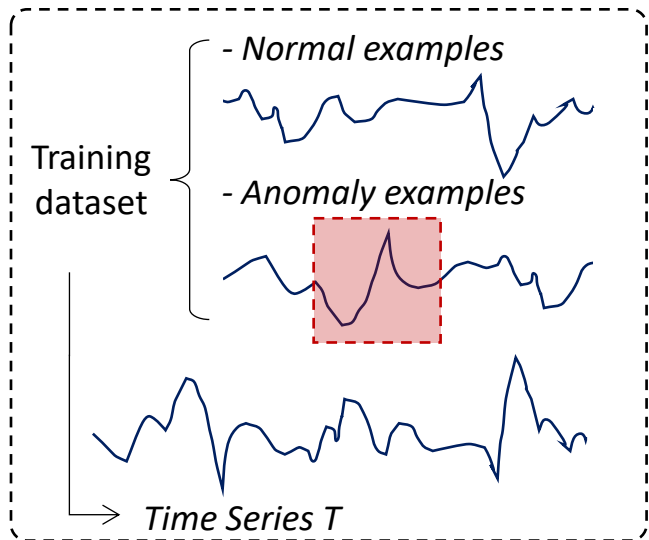
ARE

AHP

*Feed-water Pumps*

Condensor

*Cold water*

**Class 1:** Time series

Class 2

# Anomaly Dete...

By inputs...
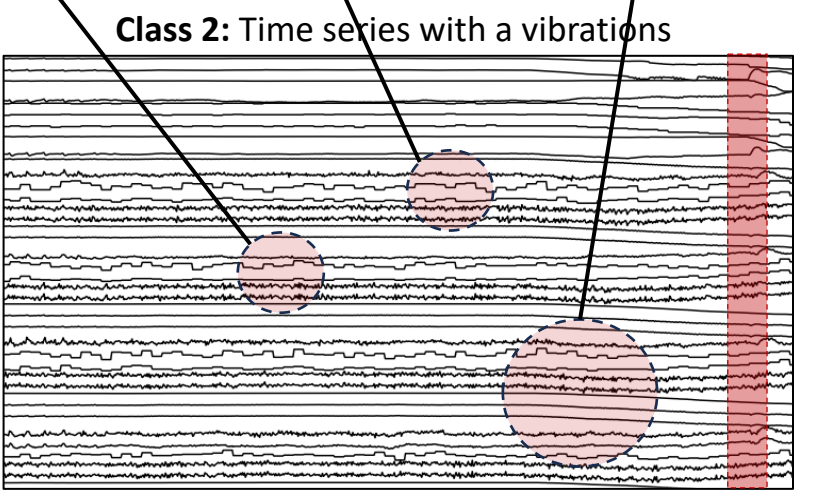


**Supervised anomaly detection (e.g., classification)**

**Explanation of the detection**

**Supervised**

**VVP: main steam system**
*28 sensors*
*(flow, pressure, temperature)*

**CEX: condensate extraction system**
*2 sensors (pressure and temperature)*

**GSS: moisture separator-reheater system**
*2 sensors (pressure and temperature)*

**GCT: turbine bypass system**
*2 sensors (pressure)*

**KKO: energy metering system**
*1 sensor (power)*

*Primary circuit*

*Secondary circuit*

CEX

VVP

*Steam*

GCT

GRE

Turbine

KKO

GSS

ASG

ADG

*High pressure water*

Steam generator

ARE

AHP

*Feed-water Pumps*

Condensor

*Cold water*

**Class 1:** Time series

**Class 2**

# Anomaly Dete...

By inputs…

**Supervised anomaly detection (e.g., classification)**

VVP: main steam system
28 sensors *(flow, pressure, temperature)*

CEX: condensate extraction system
2 sensors *(pressure and temperature)*

GSS: moisture separator-reheater system
2 sensors *(pressure and temperature)*

GCT: turbine bypass system
2 sensors *(pressure)*

KKO: energy metering system
1 sensor *(power)*

*Primary circuit*

*Secondary circuit*

VVP — *Steam*

CEX

GSS

GCT

GRE — Turbine

KKO

*(c)* **APP011MD**: *Water flow exiting the pump*

*(c.1.1)*

*(c.2.1)*

*(c.3.1)*

Cluster 7, nb element: 25

Cluster 8, nb element: 5

Cluster 9, nb element: 15

Cluster 10, nb element: 33

time_vibration - time

AGR606MT-: number of sequence: 240

number of sequence: 24

# Anomaly Dete...

By inputs...

Supervised anomaly detection (e.g., classification)

VVP: main steam system
28 sensors
(flow, pressure, temperature)

CEX: condensate extraction system
2 sensors (pressure and temperature)

GSS: moisture separator-reheater system
2 sensors (pressure and temperature)

GCT: turbine bypass system
2 sensors (pressure)

KKO: energy metering system
1 sensor (power)

Primary circuit

Secondary circuit

CEX

Steam

VVP

GRE
Turbine

KKO

pump

## More info :

### On the use case



*DCE journal 2023*

### On the method



*SIGMOD 2022*

Cluster 7, nb element: 25

Cluster 8, nb element: 5

Cluster 9, nb element: 15

# Anomaly Detection methods: *A taxonomy*

By methods…



Time series anomaly detection methods

Distance-based

Proximity-based

Clustering-based

Discord-based

*E.g.,*
LOF

*E.g.,*
NormA
SAND

*E.g.,*
MP
DAMP

# Anomaly Detection methods: *A taxonomy*

By methods…

# Anomaly Detection methods: *A taxonomy*

By methods...



Time series anomaly detection methods

Distance-based | Density-based | Prediction-based

Proximity-based | Clustering-based | Discord-based | Distribution-based | Graph-based | Tree-based | Encoding-based | Forecasting-based | Reconstruction-based

E.g., LOF | E.g., NormA SAND | E.g., MP DAMP | E.g., HOBS OCSVM | E.g., Series2Graph | E.g., Isolation-Forest | E.g., GrammarViz POLY, PCA | E.g., LSTM,CNN | E.g., AutoEncoder

# Anomaly Detection methods:
*A taxonomy*

By time…

# Anomaly Detection methods: *A taxonomy*

By time…



Number of methods proposed per Second-level categories

Number of methods proposed per Second-level categories (cumulative)

# Anomaly Detection methods: *A taxonomy*

By time…



Number of methods proposed that are
*Unsupervised* or *Semi-Supervised*

Number of methods proposed that can handle
*Univariate* or *Multivariate* time series

# Anomaly Detection methods: *Distance-based*

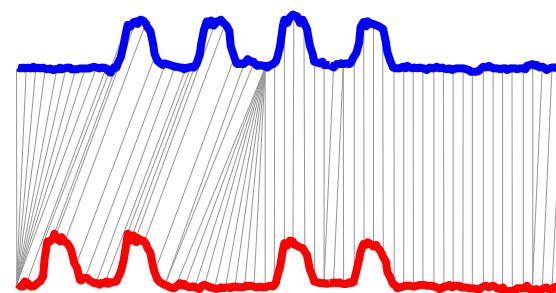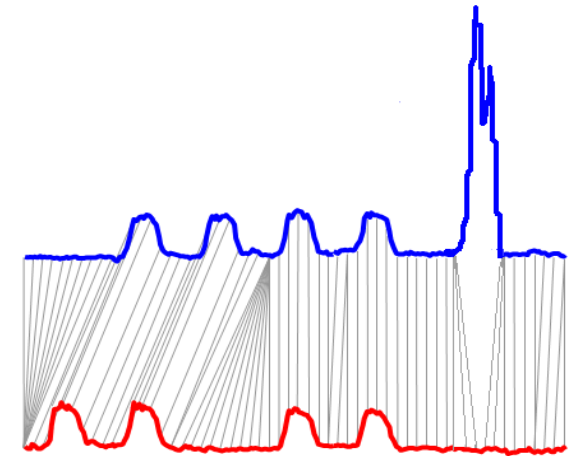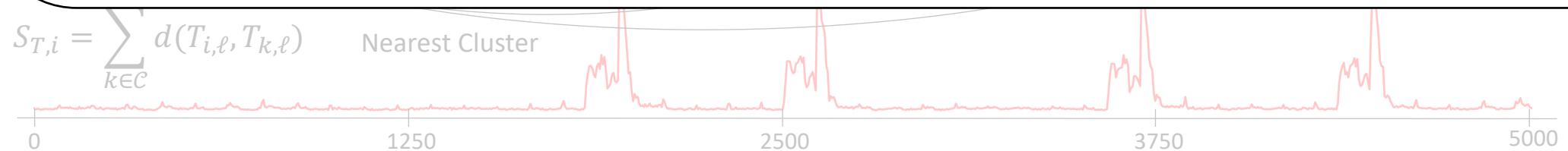Methods that use distance computation between subsequences (or group of subsequences) to detect anomalies.



Time series $T$

$T_{i,\ell}$   $T_{j,\ell}$   $T_{k,\ell}$   $T_{l,\ell}$   $T_{m,\ell}$   ...   $T_{n,\ell}$   $T_{o,\ell}$

# Anomaly Detection methods: *Distance-based*

Methods that use distance computation between subsequences (or group of subsequences) to detect anomalies.



$$S_{T,i} = d(T_{i,\ell}, T_{j,\ell})$$

# Anomaly Detection methods: *Distance-based*

Methods that use distance computation between subsequences (or group of subsequences) to detect anomalies.



$$S_{T,i} = d(T_{i,\ell}, T_{m,\ell})$$

K-Nearest neighbor

# Anomaly Detection methods: *Distance-based*

Methods that use distance computation between subsequences (or group of subsequences) to detect anomalies.

Time series $T$



$$S_{T,i} = \sum_{k \in \mathcal{C}} d(T_{i,\ell}, T_{k,\ell})$$

Nearest Cluster

# Anomaly Detection methods: *Distance-based*

Methods that use distance computation between subsequences (or group of subsequences) to detect anomalies.



$$S_{T,i} = \sum_{k \in \mathcal{C}} d(T_{i,\ell}, T_{k,\ell})$$

Nearest Cluster

# Anomaly Detection methods: *Distance-based*

Methods that use distance computation between subsequences (or group of subsequences) to detect



*Example of distance computation*

(a) Euclidean Distance

(b) DTW Distance

(c) LCSS Distance

$$S_{T,i} = \sum_{k \in \mathcal{C}} d(T_{i,\ell}, T_{k,\ell})$$

Nearest Cluster

0          1250          2500          3750          5000

# Anomaly Detection methods: *an Example*



**Matrix Profile [6] (MP)**

Compute the distance to the nearest neighbor (using the MASS algorithm z-norm Euclidean distance computation) and use it as anomaly score

Unsupervised

Univariate

sequence

[6] Chin-Chia Michael Yeh, Yan Zhu, Liudmila Ulanova, Nurjahan Begum, Yifei Ding, Hoang Anh Dau, Diego Furtado Silva, Abdullah Mueen, and Eamonn J. Keogh. 2016. Matrix Prole I: All Pairs Similarity Joins for Time Series. In ICDM.

# Anomaly Detection methods: *an Example*



$$S_T = \left[ NN(T_{0,\ell}), NN(T_{1,\ell}), \dots, NN(T_{|T|-\ell,\ell}) \right]$$

The matrix Profile is computed as follows:

Labels in figure: $NN(T_{k,\ell})$, $T_{k,\ell}$, $T_{i,\ell}$, $NN(T_{j,\ell})$, $T_{j,\ell}$, $NN(T_{i,\ell})$

## Matrix Profile [6] (MP)

Compute the distance to the nearest neighbor (using the MASS algorithm z-norm Euclidean distance computation) and use it as anomaly score

Unsupervised

Univariate

sequence

[6] Chin-Chia Michael Yeh, Yan Zhu, Liudmila Ulanova, Nurjahan Begum, Yifei Ding, Hoang Anh Dau, Diego Furtado Silva, Abdullah Mueen, and Eamonn J. Keogh. 2016. Matrix Prole I: All Pairs Similarity Joins for Time Series. In ICDM.

# Anomaly Detection methods: *an Example*

## Matrix Profile [6] (MP)

Compute the distance to the nearest neighbor (using the MASS algorithm z-norm Euclidean distance computation) and use it as anomaly score

Unsupervised

Univariate

sequence

Time series $T$

Discord

Anomaly score $S_T$

$i$

[6] Chin-Chia Michael Yeh, Yan Zhu, Liudmila Ulanova, Nurjahan Begum, Yifei Ding, Hoang Anh Dau, Diego Furtado Silva, Abdullah Mueen, and Eamonn J. Keogh. 2016. Matrix Prole I: All Pairs Similarity Joins for Time Series. In ICDM.

# Anomaly Detection methods: *an Example*

Time series $T$

Discord

## Many different extensions…

- For streaming time series: **STAMPi** [6], **DAMP** [8]
- For similar recurrent anomalies: **left-STAMP** [6]
- Anytime or ordered: **STAMP** [6], **STOMP** [7]
- For multivariate time series: **mSTAMP** [9]

0    100    200    $i$    300    400

## Matrix Profile [6] (MP)

Compute the distance to the nearest neighbor (using the MASS algorithm z-norm Euclidean distance computation) and use it as anomaly score

Unsupervised

Univariate

sequence

[6] Chin-Chia Michael Yeh, Yan Zhu, Liudmila Ulanova, Nurjahan Begum, Yifei Ding, Hoang Anh Dau, Diego Furtado Silva, Abdullah Mueen, and Eamonn J. Keogh. 2016. Matrix Prole I: All Pairs Similarity Joins for Time Series. In ICDM.

# Anomaly Detection methods: *an Example*

*Time series T*



## NormA [10]

Distance-based approach that
**summarize** the time series into
**a weighted set of subsequences**
and use the distance to them as
anomaly score

Unsupervised

Univariate

sequence

[10] Paul Boniol, Michele Linardi, Federico Roncallo, Themis Palpanas, Mohammed Meftah, and Emmanuel Remy. 2021. Unsupervised and scalable subsequence anomaly detection in large data series. The VLDB Journal 30, 6 (Nov 2021), 909–931.

# Anomaly Detection methods: *an Example*



$N_M$

$(N^0_M, w^0)$

$(N^1_M, w^1)$

$(N^n_M, w^n)$

Time series T

**NormA  [10]**

Distance-based approach that
summarize the time series into
a weighted set of subsequences
and use the distance to them as
anomaly score

Unsupervised

Univariate

sequence

[10] Paul Boniol, Michele Linardi, Federico Roncallo, Themis Palpanas, Mohammed Meftah, and Emmanuel Remy. 2021. Unsupervised and scalable subsequence anomaly detection in large data series. The VLDB Journal 30, 6 (Nov 2021), 909–931.

# Anomaly Detection methods: *an Example*



$N_M$

$(N^0{}_M, w^0)$

$(N^1{}_M, w^1)$

$(N^n{}_M, w^n)$

Time series $T$

$$for\ T_{j,\ell}\ in\ T:$$
$$d = \sum_{N^i{}_M} w^i * min_{x \in [0, \ell_{N_M} - \ell]} \{dist(T_{j,\ell}, N^i{}_{M_{x,l}})\}$$

Anomaly score $S_T$

## NormA [10]

Distance-based approach that
**summarize** the time series into
**a weighted set of subsequences**
and use the distance to them as
anomaly score

Unsupervised

Univariate

sequence

[10] Paul Boniol, Michele Linardi, Federico Roncallo, Themis Palpanas, Mohammed Meftah, and Emmanuel Remy. 2021. Unsupervised and scalable subsequence anomaly detection in large data series. The VLDB Journal 30, 6 (Nov 2021), 909–931.

# Anomaly Detection methods: *an Example*



NormA  [10]

Distance-based approach that summarize the time series into a weighted set of subsequences and use the distance to them as anomaly score

Unsupervised

Univariate

sequence

$N_M$

$(N^0{}_M, w^0)$

$(N^1{}_M, w^1)$

$(N^n{}_M, w^n)$

Time series T

$T_{j,\ell}$

$T'_{j,\ell}$

$$for\ T_{j,\ell}\ in\ T:$$

$$d = \sum_{N^i{}_M} w^i * min_{x \in [0, \ell_{N_M} - \ell]} \{ dist(T_{j,\ell}, N^i{}_{M_{x,l}}) \}$$

*Anomaly score* $S_T$

[10] Paul Boniol, Michele Linardi, Federico Roncallo, Themis Palpanas, Mohammed Meftah, and Emmanuel Remy. 2021. Unsupervised and scalable subsequence anomaly detection in large data series. The VLDB Journal 30, 6 (Nov 2021), 909–931.

# Anomaly Detection methods: *an Example*



NormA [10]

Distance-based approach that summarize the time series into a weighted set of subsequences and use the distance to them as anomaly score

Unsupervised

Univariate

sequence

SAND [25]

Distance-based approach that summarize the time series into a weighted set of subsequences, and can be updated incrementally for new arriving batches of data points

[25] Paul Boniol, John Paparrizos, Themis Palpanas, and Michael J. Franklin. 2021. SAND: streaming subsequence anomaly detection.
Proc. VLDB Endow. 14, 10 (June 2021), 1717–1729.

# Anomaly Detection methods: *Density-based*

Methods that estimate the density of the space (points or subsequences) and identify as anomalies points (or sequences)that are in low-density subspace.

Time series $T$

# Anomaly Detection methods: *Density-based*

Methods that estimate the density of the space (points or subsequences) and identify as anomalies points (or sequences)that are in low-density subspace.

Time series $T$



Tree-based approaches [11]

Distribution-based Approaches [12]

Graph-based approaches [13]

...

# Anomaly Detection methods: *an Example*

0 splits

0 splits

## Isolation Forest [11]

Density-based approach that split the space randomly and using the depth of the trees to identify anomalies

Unsupervised

Univariate/Multivariate

Point/sequence

[11] F. T. Liu, K. M. Ting and Z. -H. Zhou, "Isolation Forest," 2008 Eighth IEEE International Conference on Data Mining, Pisa, Italy, 2008, pp. 413-422

# Anomaly Detection methods: *an Example*



1 splits

1 splits

### Isolation Forest [11]

Density-based approach that split the space randomly and using the depth of the trees to identify anomalies

Unsupervised

Univariate/Multivariate

Point/sequence

[11] F. T. Liu, K. M. Ting and Z. -H. Zhou, "Isolation Forest," 2008 Eighth IEEE International Conference on Data Mining, Pisa, Italy, 2008, pp. 413-422

# Anomaly Detection methods: *an Example*

2 splits

2 splits



## Isolation Forest [11]

Density-based approach that split the space randomly and using the depth of the trees to identify anomalies

Unsupervised

Univariate/Multivariate

Point/sequence

[11] F. T. Liu, K. M. Ting and Z. -H. Zhou, "Isolation Forest," 2008 Eighth IEEE International Conference on Data Mining, Pisa, Italy, 2008, pp. 413-422

# Anomaly Detection methods: *an Example*

3 splits

3 splits



## Isolation Forest [11]

Density-based approach that split the space randomly and using the depth of the trees to identify anomalies

Unsupervised

Univariate/Multivariate

Point/sequence

[11] F. T. Liu, K. M. Ting and Z. -H. Zhou, "Isolation Forest," 2008 Eighth IEEE International Conference on Data Mining, Pisa, Italy, 2008, pp. 413-422

# Anomaly Detection methods: *an Example*

4 splits

3 splits

Isolation Forest [11]

Density-based approach that split the space randomly and using the depth of the trees to identify anomalies

Unsupervised

Univariate/Multivariate

Point/sequence

[11] F. T. Liu, K. M. Ting and Z. -H. Zhou, "Isolation Forest," 2008 Eighth IEEE International Conference on Data Mining, Pisa, Italy, 2008, pp. 413-422

# Anomaly Detection methods: *an Example*

5 splits

3 splits



## Isolation Forest [11]

Density-based approach that split the space randomly and using the depth of the trees to identify anomalies

**Unsupervised**

**Univariate/Multivariate**

**Point/sequence**

[11] F. T. Liu, K. M. Ting and Z. -H. Zhou, "Isolation Forest," 2008 Eighth IEEE International Conference on Data Mining, Pisa, Italy, 2008, pp. 413-422

# Anomaly Detection methods: *an Example*



ITree$_1$

Instance A

Instance N

## Isolation Forest [11]

Density-based approach that split the space randomly and using the depth of the trees to identify anomalies

Unsupervised

Univariate/Multivariate

Point/sequence

[11] F. T. Liu, K. M. Ting and Z. -H. Zhou, "Isolation Forest," 2008 Eighth IEEE International Conference on Data Mining, Pisa, Italy, 2008, pp. 413-422

# Anomaly Detection methods: *an Example*



**Isolation Forest [11]**

Density-based approach that split the space randomly and using the depth of the trees to identify anomalies

Unsupervised

Univariate/Multivariate

Point/sequence

[11] F. T. Liu, K. M. Ting and Z. -H. Zhou, "Isolation Forest," 2008 Eighth IEEE International Conference on Data Mining, Pisa, Italy, 2008, pp. 413-422

# Anomaly Detection methods: *an Example*



Each **node** is an ensemble of similar subsequences.

Each **edge** is associated to a weight $w$ that corresponds to the number of times a subsequence move from one node to another.

## Series2Graph [13]

Density-based approach that convert the time series into a graph and detect unusual trajectories

### Unsupervised

### Univariate

### subsequence

For a given subsequence $T_{i,\ell}$ and its corresponding path $P_{th} = < N^{(i)}, N^{(i+1)}, \dots, N^{(i+\ell)} >$, we define the normality score as follows:

$$Norm(P_{th}) = \sum_{j=i}^{i+\ell-1} \frac{w\left(N^{(j)}, N^{(j+1)}\right) \deg(N^{(j)} - 1)}{\ell}$$

[13] Paul Boniol and Themis Palpanas. 2020. Series2Graph: graph-based subsequence anomaly detection for time series. Proc. VLDB Endow. 13, 12 (August 2020), 1821–1834

# Anomaly Detection methods: *an Example*



Series2Graph [13]

Density-based approach that convert the time series into a graph and detect unusual trajectories

Unsupervised

Univariate

subsequence

[13] Paul Boniol and Themis Palpanas. 2020. Series2Graph: graph-based subsequence anomaly detection for time series. Proc. VLDB Endow. 13, 12 (August 2020), 1821–1834

# Anomaly Detection methods: *an Example*



(1) subsequence embedding

Series2Graph [13]

Density-based approach that convert the time series into a graph and detect unusual trajectories

Unsupervised

Univariate

subsequence

[13] Paul Boniol and Themis Palpanas. 2020. Series2Graph: graph-based subsequence anomaly detection for time series. Proc. VLDB Endow. 13, 12 (August 2020), 1821–1834

# Anomaly Detection methods: *an Example*



(1) subsequence embedding

(2) node creation

**Series2Graph [13]**

Density-based approach that convert the time series into a graph and detect unusual trajectories

Unsupervised

Univariate

subsequence

[13] Paul Boniol and Themis Palpanas. 2020. Series2Graph: graph-based subsequence anomaly detection for time series. Proc. VLDB Endow. 13, 12 (August 2020), 1821–1834

# Anomaly Detection methods: *an Example*



(1) subsequence embedding

(3) edge creation

### Series2Graph [13]

Density-based approach that convert the time series into a graph and detect unusual trajectories

Unsupervised

Univariate

subsequence

[13] Paul Boniol and Themis Palpanas. 2020. Series2Graph: graph-based subsequence anomaly detection for time series. Proc. VLDB Endow. 13, 12 (August 2020), 1821–1834

# Anomaly Detection methods: *an Example*



$\overrightarrow{v_{ref}}$

$max($

$\overrightarrow{r_z}$

$min(T) * \mathbf{1}_{\ell-\lambda}$

$T_2$

$\overrightarrow{v_{ref}}$

$\overrightarrow{r_z}$

$w_1 = n_1$   $N_\psi^0$

$G_{\ell_G}(\mathcal{N}, \mathcal{E})$

(1) subsequence embedding          (3) edge creation

## DADS  [26]

**Distributed** version of Series2Graph

## Series2Graph [13]

Density-based approach that **convert** the time series into a **graph** and detect **unusual trajectories**

Unsupervised

Univariate

subsequence

[26] Schneider, J., Wenig, P. & Papenbrock, T. Distributed detection of sequential anomalies in univariate time series. The VLDB Journal **30**, 579–602 (2021).

# Anomaly Detection methods: *an Example*



Snippet of SED time series [14]

Pattern following a recurrent path in the graph

Pattern following an unusual path in the graph

## Series2Graph [13]

Density-based approach that convert the time series into a graph and detect unusual trajectories

Unsupervised

Univariate

subsequence

[13] Paul Boniol and Themis Palpanas. 2020. Series2Graph: graph-based subsequence anomaly detection for time series. Proc. VLDB Endow. 13, 12 (August 2020), 1821–1834

# Anomaly Detection methods: *Forecasting-based*

Methods that aims to predict the next points based on the previous ones. The prediction error is used to detect if there is an anomaly or not.

# Anomaly Detection methods: *Forecasting-based*

Methods that aims to predict the next points based on the previous ones. The prediction error is used to detect if there is an anomaly or not.

# Anomaly Detection methods: *Forecasting-based*

Methods that aims to predict the next points based on the previous ones. The prediction error is used to detect if there is an anomaly or not.

# Anomaly Detection methods: *Forecasting-based*

Methods that aims to predict the next points based on the previous ones. The prediction error is used to detect if there is an anomaly or not.

# Anomaly Detection methods: *Forecasting-based*

Methods that aims to predict the next points based on the previous ones. The prediction error is used to detect if there is an anomaly or not.

# Anomaly Detection methods: *Forecasting-based*

Methods that aims to predict the next points based on the previous ones. The prediction error is used to detect if there is an anomaly or not.



$$S_T = \left| T_i - f(T_{i-\ell,\ell}) \right|$$

# Anomaly Detection methods: *an Example*



LSTM-AD [15]

Model that stack multiple LSTM cell and use the output to predict the next value

Semi-supervised

Univariate/Multivariate

Point/sequence

[15] Pankaj Malhotra, Lovekesh Vig, Gautam Shro, and Puneet Agarwal. 2015. Long Short Term Memory Networks for Anomaly Detection in Time Series. (2015).

# Anomaly Detection methods: *an Example*



DeepAnT [16] (CNN)

Convolutional-based approach (2 convolutional layers) taking as input a sequence and aims to predict the next value.

Semi-supervised

Univariate/Multivariate

Point/sequence

[16] M. Munir, S. A. Siddiqui, A. Dengel, and S. Ahmed. 2019. DeepAnT: A Deep Learning Approach for Unsupervised Anomaly Detection in Time Series. IEEE Access 7 (2019), 1991–2005.

# Anomaly Detection methods: *an Example*



DeepAnT [16] (CNN)

Convolutional-based approach (2 convolutional layers) taking as input a sequence and aims to predict the next value.

Semi-supervised

Univariate/Multivariate

Point/sequence

[16] M. Munir, S. A. Siddiqui, A. Dengel, and S. Ahmed. 2019. DeepAnT: A Deep Learning Approach for Unsupervised Anomaly Detection in Time Series. IEEE Access 7 (2019), 1991–2005.

# Anomaly Detection methods: *an Example*



DeepAnT [16] (CNN)

Convolutional-based approach (2 convolutional layers) taking as input a sequence and aims to predict the next value.

Semi-supervised

Univariate/Multivariate

Point/sequence

[16] M. Munir, S. A. Siddiqui, A. Dengel, and S. Ahmed. 2019. DeepAnT: A Deep Learning Approach for Unsupervised Anomaly Detection in Time Series. IEEE Access 7 (2019), 1991–2005.

# Anomaly Detection methods: *Reconstruction-based*

Methods that aims to reconstruct the time series $T$ and use the reconstruction error to detect if the time series is an anomaly or not.

Time series $T$



$f(T_{i,\ell}) = T_{i,\ell}'$

$S_{T,i} = \left\| T_{i,\ell} - T_{i,\ell}' \right\|$

# Anomaly Detection methods: *Reconstruction-based*

Methods that aims to reconstruct the time series $T$ and use the reconstruction error to detect if the time series is an anomaly or not.



Time series $T$

$f(T_{i,\ell}) = T_{i,\ell}'$

$f(T_{j,\ell}) = T_{j,\ell}'$

$S_{T,i} = \left\| T_{i,\ell} - T_{i,\ell}' \right\|$

$S_{T,j} = \left\| T_{j,\ell} - T_{j,\ell}' \right\|$

# Anomaly Detection methods: *Reconstruction-based*

Methods that aims to reconstruct the time series $T$ and use the reconstruction error to detect if the time series is an anomaly or not.

Time series $T$



$f(T_{i,\ell}) = T_{i,\ell}'$

$f(T_{j,\ell}) = T_{j,\ell}'$

$S_{T,i} = \left\| T_{i,\ell} - T_{i,\ell}' \right\|$

$S_{T,j} = \left\| T_{j,\ell} - T_{j,\ell}' \right\|$

# Anomaly Detection methods: *an Example*



$$\text{Anomaly score}$$
$$S = \mathcal{L}(T_{i,l}, T'_{i,l})$$

**AutoEncoders [17] (AE)**

Neural Network composed of an encoder (that reduce the dimensionality) and decoder that reconstruct the time series. The objective is to minimize the reconstruction error.

Semi-supervised

Univariate/Multivariate

Point/sequence

[17] Mayu Sakurada and Takehisa Yairi. 2014. Anomaly Detection Using Autoencoders with Nonlinear Dimensionality Reduction. In Proceedings of the MLSDA 2014 2nd Workshop on Machine Learning for Sensory Data Analysis (Gold Coast, Australia QLD, Australia) (MLSDA'14).

# Anomaly Detection methods: *Existing benchmark*

# Anomaly Detection methods: *Existing benchmark*

## HEX/UCR [18]

Set of 250 time series with labels.

Details

- The labels have been manually checked and are reliable

- Each time series contains only 1 labeled anomaly

## TimeEval [5]

Set of 976 time series with labels.

Details

- New synthetic benchmark GutenTag used to tune parameters

- Only Time series with low contamination rate (< 0.1)

- Time series with at least one methods above 0.8 AUC-ROC

## TSB-UAD [19]

Set of 2000 time series with labels.

Details

- Collected as proposed in the literature (no filtering based on contamination, size or label quality)

- Artificial and synthetic data generation methods for reliable labels

# Anomaly Detection methods: *Existing benchmark*

# Anomaly Detection methods:
## *Experimental evaluation*

**Observations on TimeEval [5]:**

- Distance-based and Density-based methods have a better accuracy (AUC-ROC) than forecasting and reconstruction-based approaches

[5] Sebastian Schmidl, Phillip Wenig, and Thorsten Papenbrock. 2022. Anomaly detection in time series: a comprehensive evaluation. Proc. VLDB Endow. 15, 9 (May 2022), 1779–1797.



| Methods | | | | AUC-ROC |
|---|---|---|---|---|
| Sub-LOF [22] | 2 % | 0 % | 0 % | |
| GrammarViz [120] | 3 % | 0 % | 0 % | |
| DWT-MLEAD [134] | 0 % | 0 % | 0 % | |
| VALMOD [82] | 1 % | 9 % | 11 % | |
| SAND [17] | 5 % | 1 % | 22 % | |
| Left STAMPi [156] | 2 % | 0 % | 1 % | |
| Series2Graph [16] | 0 % | 0 % | 5 % | |
| ARIMA [65] | 7 % | 0 % | 0 % | |
| PCI [157] | 0 % | 0 % | 0 % | |
| STOMP [164] | 2 % | 0 % | 0 % | |
| STAMP [156] | 4 % | 0 % | 0 % | |
| Triple ES [1] | 15 % | 0 % | 9 % | |
| NumentaHTM [3] | 0 % | 0 % | 0 % | |
| NormA-SJ [15] | 10 % | 1 % | 3 % | |
| Sub-IF [83] | 0 % | 0 % | 0 % | |
| MedianMethod [10] | 0 % | 0 % | 0 % | |
| SR [112] | 0 % | 0 % | 0 % | |
| PS-SVM [85] | 12 % | 0 % | 0 % | |
| PST [128] | 0 % | 4 % | 0 % | |
| SSA [155] | 1 % | 0 % | 0 % | |
| HOT SAX [70] | 24 % | 1 % | 1 % | |
| TSBitmap [144] | 0 % | 0 % | 0 % | |
| DSPOT [122] | 6 % | 0 % | 0 % | |
| FFT [111] | 0 % | 0 % | 0 % | |
| S-H-ESD [62] | 0 % | 0 % | 49 % | |
| Donut [150] | 1 % | 1 % | 2 % | |
| RForest [21] | 12 % | 0 % | 0 % | |
| IE-CAE [44] | 0 % | 0 % | 1 % | |
| XGBoosting [34] | 0 % | 0 % | 0 % | |
| OceanWNN [143] | 0 % | 0 % | 10 % | |
| Bagel [79] | 19 % | 0 % | 2 % | |
| SR-CNN [112] | 22 % | 0 % | 1 % | |
| TARZAN [71] | 0 % | 0 % | 18 % | |

(Rows Sub-LOF through S-H-ESD: **Unsupervised**; rows Donut through TARZAN: **Semi-supervised**)

# Anomaly Detection methods:
## *Experimental evaluation*

**Observations on TimeEval [5]:**

- Distance-based and Density-based methods have a better accuracy (AUC-ROC) than forecasting and reconstruction-based approaches

- Semi-supervised methods are not outperforming Unsupervised approaches

[5] Sebastian Schmidl, Phillip Wenig, and Thorsten Papenbrock. 2022. Anomaly detection in time series: a comprehensive evaluation. Proc. VLDB Endow. 15, 9 (May 2022), 1779–1797.



| Methods | | | | AUC-ROC |
|---|---|---|---|---|
| **Unsupervised** | | | | |
| Sub-LOF [22] | 2 % | 0 % | 0 % | |
| GrammarViz [120] | 3 % | 0 % | 0 % | |
| DWT-MLEAD [134] | 0 % | 0 % | 0 % | |
| VALMOD [82] | 1 % | 9 % | 11 % | |
| SAND [17] | 5 % | 1 % | 22 % | |
| Left STAMPi [156] | 2 % | 0 % | 1 % | |
| Series2Graph [16] | 0 % | 0 % | 5 % | |
| ARIMA [65] | 7 % | 0 % | 0 % | |
| PCI [157] | 0 % | 0 % | 0 % | |
| STOMP [164] | 2 % | 0 % | 0 % | |
| STAMP [156] | 4 % | 0 % | 0 % | |
| Triple ES [1] | 15 % | 0 % | 9 % | |
| NumentaHTM [3] | 0 % | 0 % | 0 % | |
| NormA-SJ [15] | 10 % | 1 % | 3 % | |
| Sub-IF [83] | 0 % | 0 % | 0 % | |
| MedianMethod [10] | 0 % | 0 % | 0 % | |
| SR [112] | 0 % | 0 % | 0 % | |
| PS-SVM [85] | 12 % | 0 % | 0 % | |
| PST [128] | 0 % | 4 % | 0 % | |
| SSA [155] | 1 % | 0 % | 0 % | |
| HOT SAX [70] | 24 % | 1 % | 1 % | |
| TSBitmap [144] | 0 % | 0 % | 0 % | |
| DSPOT [122] | 6 % | 0 % | 0 % | |
| FFT [111] | 0 % | 0 % | 0 % | |
| S-H-ESD [62] | 0 % | 0 % | 49 % | |
| **Semi-supervised** | | | | |
| Donut [150] | 1 % | 1 % | 2 % | |
| RForest [21] | 12 % | 0 % | 0 % | |
| IE-CAE [44] | 0 % | 0 % | 1 % | |
| XGBoosting [34] | 0 % | 0 % | 0 % | |
| OceanWNN [143] | 0 % | 0 % | 10 % | |
| Bagel [79] | 19 % | 0 % | 2 % | |
| SR-CNN [112] | 22 % | 0 % | 1 % | |
| TARZAN [71] | 0 % | 0 % | 18 % | |

# Anomaly Detection methods:
*Experimental evaluation*

**Observations on HEX/UCR [18]:**

- Distance-based methods have a better accuracy (AUC-ROC) than forecasting and distribution-based approaches

[18] R. Wu and E. Keogh, "Current Time Series Anomaly Detection Benchmarks are Flawed and are Creating the Illusion of Progress" in IEEE Transactions on Knowledge & Data Engineering, vol. 35, no. 03, pp. 2421-2429, 2023.

# Anomaly Detection methods:
## *Experimental evaluation*

**Observations on TSB-UAD [19]:**

- Distance-based methods have a better accuracy (AUC-ROC) than forecasting-based methods.
- Isolation Forest (Tree-based and not proposed for time series) have also a strong accuracy

- AutoEncoder (AE) is also very accurate.

[19] John Paparrizos, Yuhao Kang, Paul Boniol, Ruey S. Tsay, Themis Palpanas, and Michael J. Franklin. 2022. TSB-UAD: an end-to-end benchmark suite for univariate time-series anomaly detection. Proc. VLDB Endow. 15, 8 (April 2022), 1697–1711.

# Anomaly Detection methods:
## *Experimental evaluation*

**Point-based** anomaly

**sequence-based** anomaly

**Observations on TSB-UAD [19]:**

- Forecasting methods (LSTM and CNN) are very accurate for point anomalies
- But have poor performances on sequence-based anomalies.

[19] John Paparrizos, Yuhao Kang, Paul Boniol, Ruey S. Tsay, Themis Palpanas, and Michael J. Franklin. 2022. TSB-UAD: an end-to-end benchmark suite for univariate time-series anomaly detection. Proc. VLDB Endow. 15, 8 (April 2022), 1697–1711.

# Anomaly Detection methods:
## *Experimental evaluation*

**Observations on TSB-UAD [19]:**

- The ratio of normal/abnormal points has a
  <span style="color:red">strong impact</span> on the methods ranking.

[19] John Paparrizos, Yuhao Kang, Paul Boniol, Ruey S. Tsay, Themis Palpanas, and Michael J. Franklin. 2022. TSB-UAD: an end-to-end benchmark suite for univariate time-series anomaly detection. Proc. VLDB Endow. 15, 8 (April 2022), 1697–1711.

MGAB     MITDB     Genesis     Dodgers

Daphnet     MGAB     MITDB     SMD

SensorScope     NASA     SMD     YAHOO

Genesis     Dodgers     YAHOO

[19] John Paparrizos, Yuhao Kang, Paul Boniol, Ruey S. Tsay, Themis Palpanas, and Michael J. Franklin. 2022. TSB-UAD: an end-to-end benchmark suite for univariate time-series anomaly detection. Proc. VLDB Endow. 15, 8 (April 2022), 1697–1711.

# Perspectives and challenges

# Conclusion and Open Problems

If you are interested in anomaly detection in time series…



https://github.com/HPI-Information-Systems/TimeEval

S. Schmidl et al. PVLDB (2022) [5]

https://github.com/TheDatumOrg/TSB-UAD

J. Paparrizos et al. PVLDB (2022) [19]

https://wu.renjie.im/research/anomaly-benchmarks-are-flawed/

R. Wu et al. TKDE (2021) [18]

A. Blazquez-Garcia et al. ACM Computing Survey (2021) [24]

# Conclusion and Open Problems

If you are interested in anomaly dete



TSB
*kit*

*PyPI v0.0.5, Python 3.8*

`Pip install tsb-kit`

GitHub

Documentation

**Anomaly Detection in Time Series: A Comprehensive Evaluation**

Sebastian Schmidl[*]
Hasso Plattner Institute,
University of Potsdam
Potsdam, Germany
sebastian.schmidl@hpi.de

Phillip Wenig[*]
Hasso Plattner Institute,
University of Potsdam
Potsdam, Germany
phillip.wenig@hpi.de

Thorsten Papenbrock
Philipps University of Marburg
Marburg, Germany
papenbrock@informatik.uni-marburg.de

https://github.com/HPI-Information-Systems/TimeEval

S. Schmidl et al. PVLDB (2022) [5]

**TSB-UAD: An End-t
Time-**

John Paparrizos
The Ohio State University
paparrizos.1@osu.edu

Ruey S. Tsay
University of Chicago
ruey.tsay@chicagobooth.edu

https://githu

J. Paparrizo

...maly Detection
...d are Creating the
...ress

...Keogh

...E (2021)

**A review on outlier/anomaly detection in time series data**

ANE BLÁZQUEZ-GARCÍA and ANGEL CONDE, Ikerlan Technology Research Centre, Basque Research and Technology Alliance (BRTA), Spain

USUE MORI, Intelligent Systems Group (ISG), Department of Computer Science and Artificial Intelligence, University of the Basque Country (UPV/EHU), Spain

JOSE A. LOZANO, Intelligent Systems Group (ISG), Department of Computer Science and Artificial Intelligence, University of the Basque Country (UPV/EHU), Spain and Basque Center for Applied Mathematics (BCAM), Spain

A. Blazquez-Garcia et al. ACM Computing Survey (2021) [24]

# Conclusion and Open Problems

Context-aware Unsupervised Anomaly Detection



*number of taxi passengers in New York City*

Christmas week

Daylight Saving Time (DST)

Flooding

Snowstorm

# Conclusion and Open Problems

Evaluating Anomaly Detection



(ex1) *Example on IOPS*

(ex2) *Example on SensorScope*

(ex3) *Example on NAB*

## What is the problem here?

# Conclusion and Open Problems

Threshold-based Evaluation
Measures:



Labels

Time Series

Thresholds T

Anomaly score

# Conclusion and Open Problems

Threshold-based Evaluation
Measures:

Labels

Time Series

0

8000

0

8000

Anor

# Conclusion and Open Problems

Threshold-based Evaluation Measures:

- Precision: $\frac{TP}{TP+FP}$

- Recall (true positive rate): $\frac{TP}{TP+FN}$

- False positive rate: $\frac{FP}{FP+TN}$

- F-score: $\frac{(1+\beta^2)*Precision}{\beta^2*Precision+Recall}$

# Conclusion and Open Problems

**Labeling can be an issue for time series [22]:**

- Misalignment can lead to significant changes of accuracy values.

- This is a real issue because of:

  - Methods that produce misaligned anomaly scores.

  - Different Labeling strategies between domains and applications

(1) *Time series*

anomaly

(2) *Anomaly score*
— : *Subsequence method* ($\ell$)
— : *Point method*

(3) *Labeling strategy:*
*anomaly:*
*anomaly + borders:*
*anomaly + right border:*

(ex1) *Example on IOPS*

(ex2) *Example on SensorScope*

(ex3) *Example on NAB*

# Conclusion and Open Problems

If you are interested in evaluation measures for anomaly detection...

# Conclusion and Open Problems

If you are interested in evaluation measures for anomaly detection…



**Precision and Recall for Time Series**

https://arxiv.org/abs/1803.03639

N. Tatbul et al. NeurIPS 2018 [23]



**Volume Under the Surface: A New Accuracy Evaluation Measure for Time-Series Anomaly Detection**

https://www.vldb.org/pvldb/vol15/p2774-paparrizos.pdf

J. Paparrizos et al. PVLDB 2022 [22]



**Local Evaluation of Time Series Anomaly Detection Algorithms**

https://arxiv.org/abs/2206.13167

A. Huet et al. KDD 2022 [31]



**Navigating the Metric Maze: A Taxonomy of Evaluation Metrics for Anomaly Detection in Time Series**

https://arxiv.org/abs/2303.01272

S. Sørbø et al. DAMI 2024 [29]

NAB   Genesis   Dodgers   YAHOO   MITDB

ECG   GHL   SensorScope   NASA-MSL   SMD

NASA-SMAP   NAB   Genesis   Dodgers   YAHOO

Methods ranking changes significantly between datasets [19]

Results over TSB-UAD

(a.1) Example from ECG dataset   (a.2) ECG

(b.1) Example from MGAB dataset   (b.2) MGAB

(c.1) Example from Daphnet dataset   (c.2) Daphnet

(d.1) Example from YAHOO dataset   (d.2) YAHOO

NAB  Genesis  Dodgers  YAHOO  MITDB

ECG  GHL  SensorScope  NASA-MSL  SMD

NASA-SMAP  NAB  Genesis  Dodgers  YAHOO

Methods ranking changes significantly between datasets [19]

Can *Ensembling* methods solve the problem?

(a.1) Example from ECG dataset  (a.2) ECG

(b.1) Example from MGAB dataset  (b.2) MGAB

(c.1) Example from Daphnet dataset  (c.2) Daphnet

(d.1) Example from YAHOO dataset  (d.2) YAHOO

NAB  Genesis  Dodgers  YAHOO  MITDB

ECG  GHL  SensorScope  NASA-MSL  SMD

NASA-SMAP  NAB  Genesis  Dodgers  YAHOO

Methods ranking changes significantly between datasets [19]

Can *automatic model selection solve the problem*?

(a.1) Example from ECG dataset  (a.2) ECG

(b.1) Example from MGAB dataset  (b.2) MGAB

(c.1) Example from Daphnet dataset  (c.2) Daphnet

(d.1) Example from YAHOO dataset  (d.2) YAHOO

NAB  Genesis  Dodgers  YAHOO  MITDB
ECG  GHL  SensorScope  NASA-MSL  SMD
NASA-SMAP  NAB  Genesis  Dodgers  YAHOO

Methods ranking changes significantly between datasets [19]

Can *automatic model selection solve the problem*?

(a.1) Example from ECG dataset     (a.2) ECG
(b.1) Example from MGAB dataset     (b.2) MGAB
(c.1) Example from Daphnet dataset  (c.2) Daphnet
(d.1) Example from YAHOO dataset    (d.2) YAHOO

NAB · Genesis · Dodgers · YAHOO · MITDB

ECG · GHL · SensorScope · NASA-MSL · SMD

NASA-SMAP · NAB · Genesis · Dodgers · YAHOO

Methods ranking changes significantly between datasets [19]

(a.1) Example from ECG dataset (a.2) ECG

(b.1) Example from MGAB dataset (b.2) MGAB

(c.1) Example from Daphnet dataset (c.2) Daphnet

(d.1) Example from YAHOO dataset (d.2) YAHOO

IFOREST, LOF, MP, NORMA, IFOREST1, HBOS, OCSVM, PCA, AE, CNN, LSTM, POLY

Can *automatic model selection solve the problem*?

## Choose Wisely [29]

An experimental evaluation of model selection for time series anomaly detection

*VLDB 2023*                    *ICDE 2024*

[29] Emmanouil Sylligardos, Paul Boniol, John Paparrizos, Panos Trahanias, and Themis Palpanas. 2023. Choose Wisely: An Extensive Evaluation of Model Selection for Anomaly Detection in Time Series. Proc. VLDB Endow. 16, 11 (July 2023), 3418–3432.

# References

[1] N. Laptev, S. Amizadeh, and Y. Billawala. 2015. **S5 - A Labeled Anomaly Detection Dataset**, version 1.0(16M).

[2] Markus Thill, Wolfgang Konen, and Thomas Bäck. 2020. **MGAB: The Mackey-Glass Anomaly Benchmark**.

[3] Pawel Benecki, Szymon Piechaczek, Daniel Kostrzewa, and Jakub Nalepa. 2021. **Detecting Anomalies in Spacecraft Telemetry Using Evolutionary Thresholding and LSTMs**. In Proceedings of the Genetic and Evolutionary Computation Conference Companion (Lille, France) (GECCO '21)

[4] Scott David Greenwald. 1990. **Improved detection and classication of arrhythmias in noise-corrupted electrocardiograms using contextual information**. Thesis. Massachusetts Institute of Technology.

[5] Sebastian Schmidl, Phillip Wenig, and Thorsten Papenbrock. 2022. **Anomaly detection in time series: a comprehensive evaluation**. Proc. VLDB Endow. 15, 9 (May 2022), 1779–1797.

[6] Chin-Chia Michael Yeh, Yan Zhu, Liudmila Ulanova, Nurjahan Begum, Yifei Ding, Hoang Anh Dau, Diego Furtado Silva, Abdullah Mueen, and Eamonn J. Keogh. 2016. **Matrix Prole I: All Pairs Similarity Joins for Time Series.** In ICDM.

[7] Yan Zhu, Zachary Zimmerman, Nader Shakibay Senobari, Chin-Chia Michael Yeh, Gareth Funning, Abdullah Mueen, Philip Brisk, and Eamonn Keogh. 2016. **Matrix Profile II: Exploiting a Novel Algorithm and GPUs to Break the One Hundred Million Barrier for Time Series Motifs and Joins.** In Proceedings of the International Conference on Data Mining (ICDM), 739–748.

[8] Yue Lu, Renjie Wu, Abdullah Mueen, Maria A. Zuluaga, and Eamonn Keogh. 2022. **Matrix Profile XXIV: Scaling Time Series Anomaly Detection to Trillions of Datapoints and Ultra-fast Arriving Data Streams.** In Proceedings of the 28th ACM SIGKDD.

[9] C. -C. M. Yeh, N. Kavantzas and E. Keogh, **Matrix Profile VI: Meaningful Multidimensional Motif Discovery**, 2017 IEEE International Conference on Data Mining (ICDM), New Orleans, LA, USA, 2017, pp. 565-574, doi: 10.1109/ICDM.2017.66. Data Mining (KDD '22).

[10] Paul Boniol, Michele Linardi, Federico Roncallo, Themis Palpanas, Mohammed Meftah, and Emmanuel Remy. 2021. **Unsupervised and scalable subsequence anomaly detection in large data series.** The VLDB Journal 30, 6 (Nov 2021), 909–931.

[11] F. T. Liu, K. M. Ting and Z. -H. Zhou, **Isolation Forest**, 2008 Eighth IEEE International Conference on Data Mining, Pisa, Italy, 2008, pp. 413-422

[12] Markus Goldstein and Andreas Dengel. 2012. **Histogram-based outlier score (hbos): A fast unsupervised anomaly detection algorithm**. KI-2012: poster and demo track 9 (2012).

[13] Paul Boniol and Themis Palpanas. 2020. **Series2Graph: graph-based subsequence anomaly detection for time series.** Proc. VLDB Endow. 13, 12 (August 2020), 1821–1834.

[14] Ali Abdul-Aziz, Mark R Woike, Nikunj C Oza, Bryan L Matthews, and John D lekki. 2012. **Rotor health monitoring combining spin tests and data-driven anomaly detection methods**. Structural Health Monitoring (2012).

[15] Pankaj Malhotra, Lovekesh Vig, Gautam Shro, and Puneet Agarwal. 2015. **Long Short Term Memory Networks for Anomaly Detection in Time Series**. (2015).

[16] M. Munir, S. A. Siddiqui, A. Dengel, and S. Ahmed. 2019. **DeepAnT: A Deep Learning Approach for Unsupervised Anomaly Detection in Time Series**. IEEE Access 7 (2019), 1991–2005.

[17] Mayu Sakurada and Takehisa Yairi. 2014. **Anomaly Detection Using Autoencoders with Nonlinear Dimensionality Reduction**. In Proceedings of the MLSDA 2014 2nd Workshop on Machine Learning for Sensory Data Analysis (Gold Coast, Australia QLD, Australia) (MLSDA'14).

[18] R. Wu and E. Keogh, **Current Time Series Anomaly Detection Benchmarks are Flawed and are Creating the Illusion of Progress** in IEEE Transactions on Knowledge & Data Engineering, vol. 35, no. 03, pp. 2421-2429, 2023.

[19] John Paparrizos, Yuhao Kang, Paul Boniol, Ruey S. Tsay, Themis Palpanas, and Michael J. Franklin. 2022. **TSB-UAD: an end-to-end benchmark suite for univariate time-series anomaly detection**. Proc. VLDB Endow. 15, 8 (April 2022), 1697–1711.

[20] Tom Fawcett. 2006. **An introduction to ROC analysis**. Pattern Recognition Letters 27, 8 (2006), 861–874.

[21] Jesse Davis and Mark Goadrich. 2006. **The Relationship between Precision-Recall and ROC Curves**. In Proceedings of the 23rd International Conference on Machine Learning (ICML '06).

[22] John Paparrizos, Paul Boniol, Themis Palpanas, Ruey S. Tsay, Aaron Elmore, and Michael J. Franklin. 2022. **Volume under the surface: a new accuracy evaluation measure for time-series anomaly detection**. Proc. VLDB Endow. 15, 11 (July 2022), 2774–2787.

[23] Nesime Tatbul, Tae Jun Lee, Stan Zdonik, Mejbah Alam, and Justin Gottschlich. 2018. **Precision and Recall for Time Series**. In Advances in Neural Information Processing Systems, Vol. 31.

[24] Ane Blázquez-García, Angel Conde, Usue Mori, and Jose A. Lozano. 2021. **A Review on Outlier/Anomaly Detection in Time Series Data**. ACM Comput. Surv. 54, 3, Article 56 (April 2022), 33 pages.

[25] Paul Boniol, John Paparrizos, Themis Palpanas, and Michael J. Franklin. 2021. **SAND: streaming subsequence anomaly detection**. Proc. VLDB Endow. 14, 10 (June 2021), 1717–1729.

[26] Schneider, J., Wenig, P. & Papenbrock, T. **Distributed detection of sequential anomalies in univariate time series**. The VLDB Journal 30, 579–602 (2021).

[27] Paul Boniol, John Paparrizos, Yuhao Kang, Themis Palpanas, Ruey S. Tsay, Aaron J. Elmore, and Michael J. Franklin. 2022. **Theseus: navigating the labyrinth of time-series anomaly detection.** Proc. VLDB Endow. 15, 12 (August 2022), 3702–3705.

[28] Paul Boniol, Themis Palpanas, Mohammed Meftah, and Emmanuel Remy. 2020. **GraphAn: graph-based subsequence anomaly detection.** Proc. VLDB Endow. 13, 12 (August 2020), 2941–2944.

[29] Sørbø, S., Ruocco, M. **Navigating the metric maze: a taxonomy of evaluation metrics for anomaly detection in time series**. Data Min Knowl Disc 38, 1027–1068 (2024)

[30] Emmanouil Sylligardos, Paul Boniol, John Paparrizos, Panos Trahanias, and Themis Palpanas. 2023. **Choose Wisely: An Extensive Evaluation of Model Selection for Anomaly Detection in Time Series.** Proc. VLDB Endow. 16, 11 (July 2023), 3418–3432.

[31] Alexis Huet, Jose Manuel Navarro, and Dario Rossi. 2022. **Local Evaluation of Time Series Anomaly Detection Algorithms.** In Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining (KDD '22). Association for Computing Machinery, New York, NY, USA, 635–645.

# Thank you for attending!

---

## Any Questions?