



Instituto de Comunicação e Informação  
Científica e Tecnológica em Saúde



CEFET/RJ

DAL

# OUTBREAKS DETECTION

Eduardo Ogasawara  
eogasawara@ieee.org  
<https://eic.cefet-rj.br/~eogasawara>

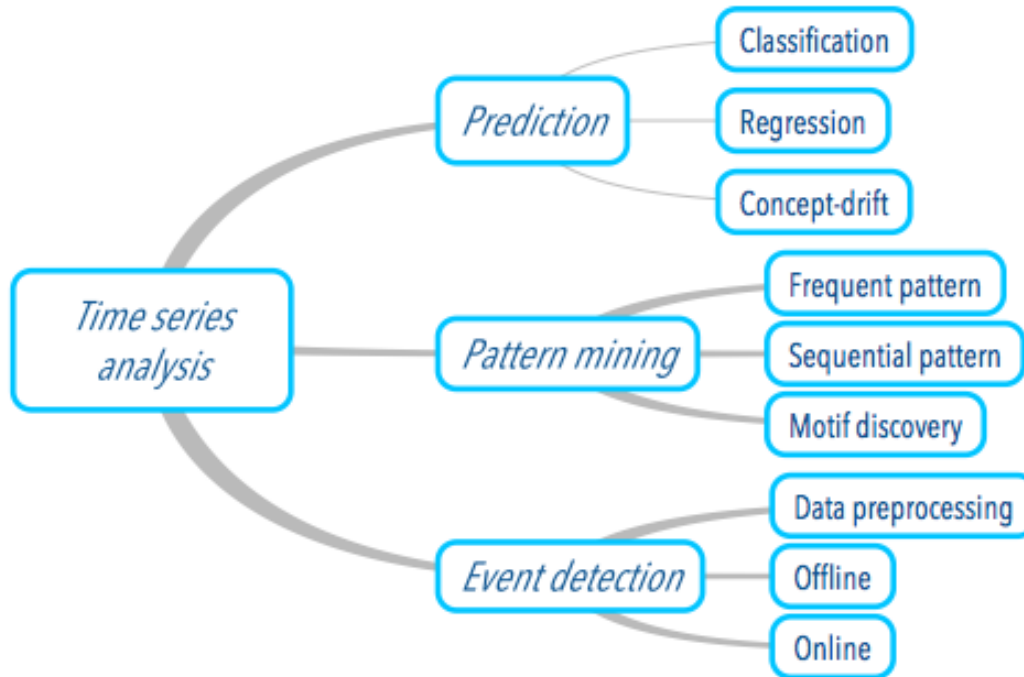
# Biograph

- D.Sc. in Systems and Computer Engineering (COPPE/UFRJ) in 2011
- Professor at EIC - CEFET/RJ
  - Computer Science Department
  - Computer Science Technical High School
- Graduate Program in Computer Science (PPCIC)
- Graduate Program in Eng. Production and Systems (PPPRO)
- Member of IEEE, SBC, ACM, and INNS
- Associated Editor of IEEE Latin America Transactions



<https://eic.cefet-rj.br/~eogasawara>

# Research Themes



Let's start

# Disease Outbreaks

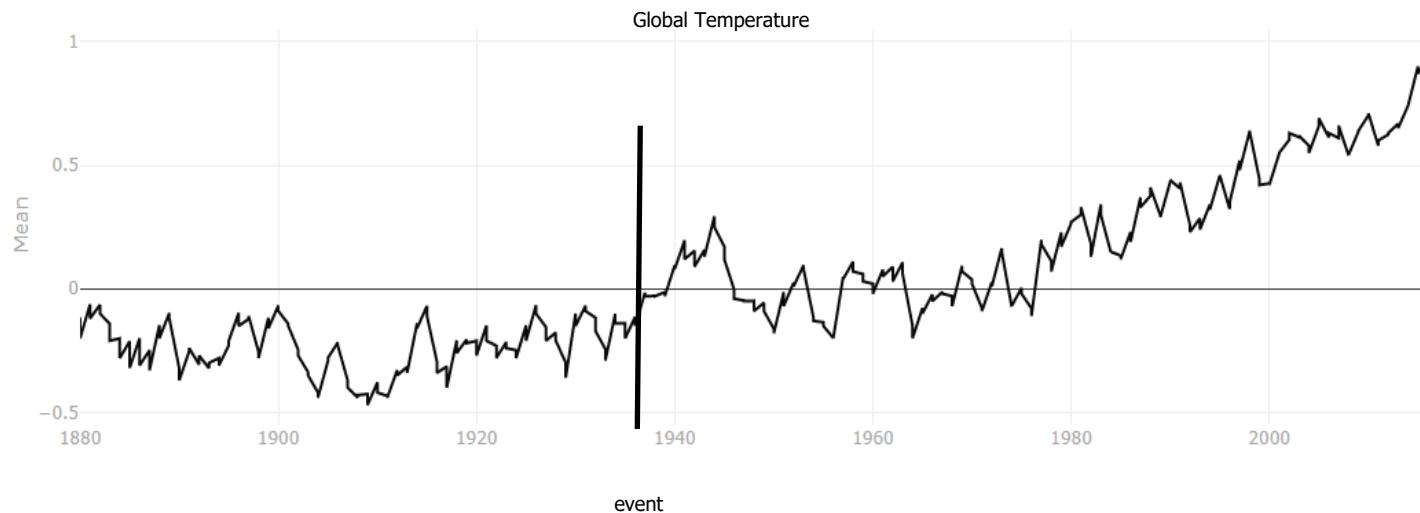
- A disease outbreak is the occurrence of disease cases in excess of normal expectancy
  - The number of cases varies according to the disease-causing agent, and the size and type of previous and existing exposure to the agent
- Disease outbreaks are usually caused by an infection, transmitted through person-to-person contact, animal-to-person contact, or from the environment or other media
  - Outbreaks may also occur following exposure to chemicals or to radioactive materials



Source: WHO,2022

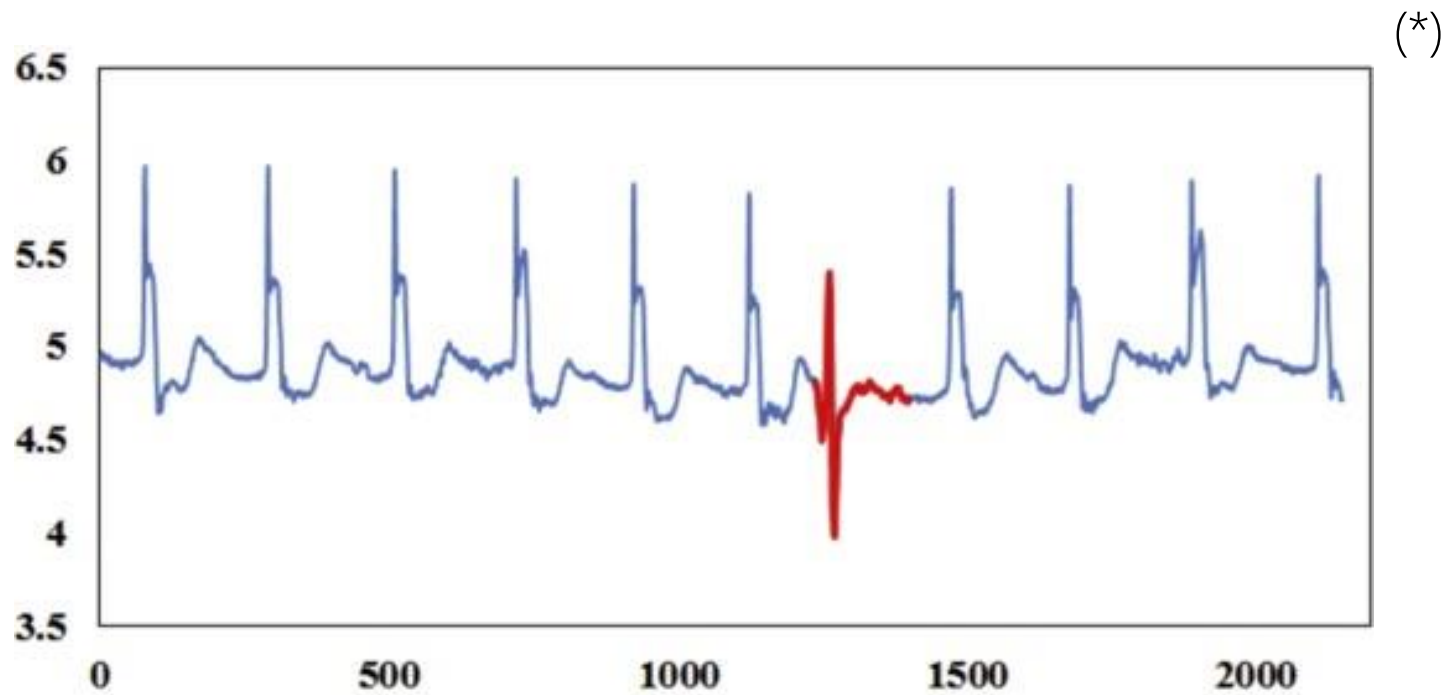
# Events

- A point or an interval where a significant change in the time series behavior occurs
- Events may appear as anomalies, change points, or frequent patterns (motifs)



# Anomalies

- A pattern or observation that do not conform to expected behavior [1]
- It can be categorized as punctual, contextual or collective

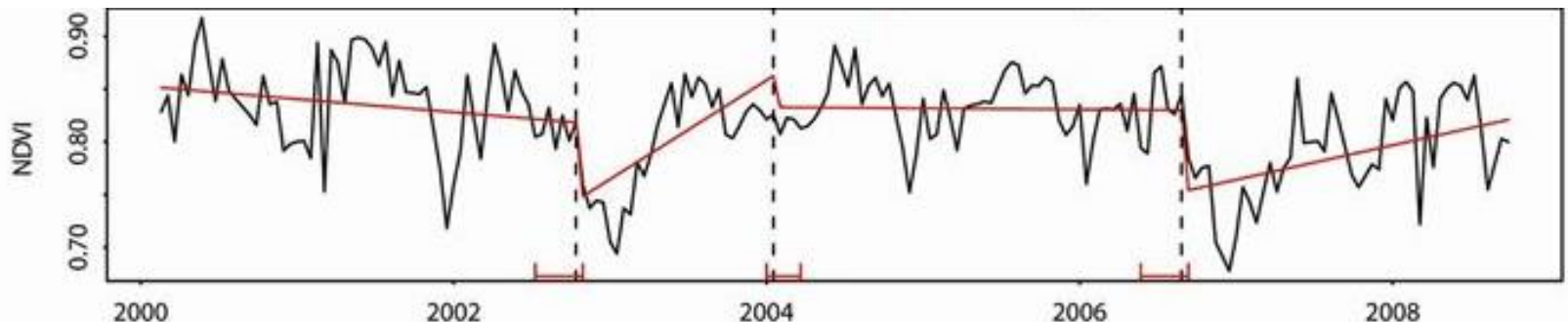


[1] V. Chandola, A. Banerjee, e V. Kumar, 2009, Anomaly detection: A survey, ACM Computing Surveys, v. 41, n. 3

(\*) In this example, it can also be classified as a discord

# Change Points

- Points (or time intervals) that mark significant change in time series behavior [1]
- They separate different states in the process that generates the time series

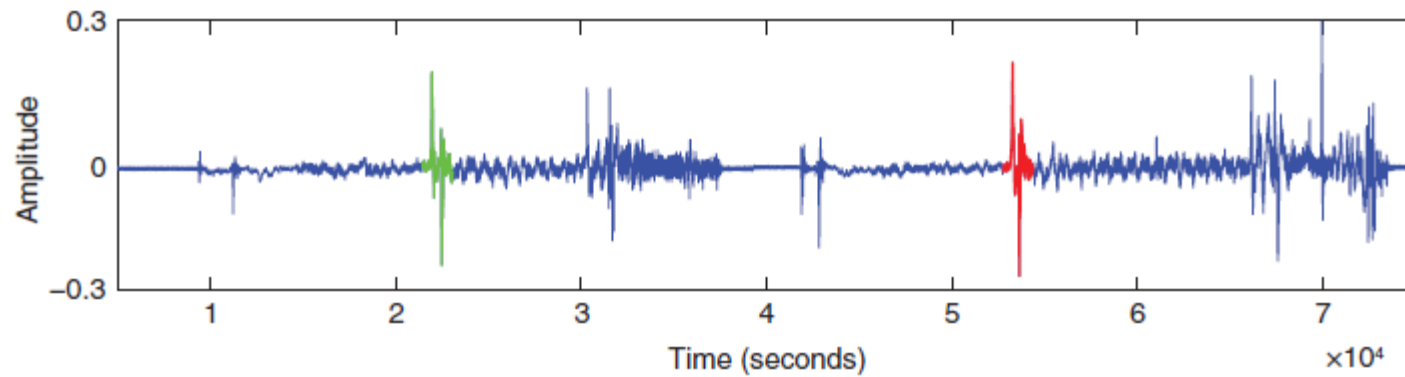


[1] J.-I. Takeuchi e K. Yamanishi, 2006, A unifying framework for detecting outliers and change points from time series, IEEE Transactions on Knowledge and Data Engineering, v. 18, n. 4, p. 482–492.



# Motifs

- A pattern (unknown) that occurs a significant number of times in time series [1,2,3]



How to do it in non-stationarity time series?

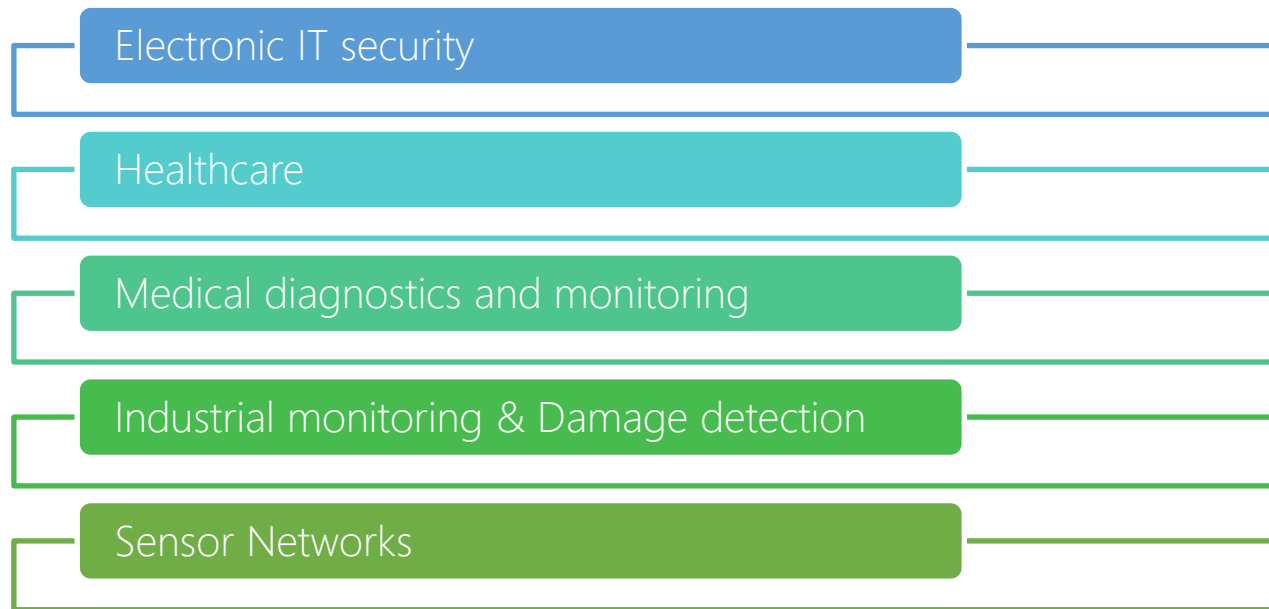
[1] P. Patel, E. Keogh, J. Lin, and S. Lonardi, "Mining motifs in massive time series databases," in Proceedings - IEEE International Conference on Data Mining, ICDM, 2002, pp. 370–377

[2] A. Mueen, "Time series motif discovery: Dimensions and applications," Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery, vol. 4, no. 2, pp. 152–159, 2014

[3] S. Torkamani and V. Lohweg, "Survey on time series motif discovery," Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery, vol. 7, no. 2, 2017.

# Event detection

- An event can represent a phenomenon with a specific meaning defined in a certain domain
- Event detection is the process of finding events
- It is a basic function in surveillance and monitoring systems
- Example of applications:



[1] V. Chandola, A. Banerjee, e V. Kumar, 2009, Anomaly detection: A survey, ACM Computing Surveys, v. 41, n. 3

[2] M. Gupta, J. Gao, C.C. Aggarwal, e J. Han, 2014, Outlier Detection for Temporal Data: A Survey, IEEE Transactions on Knowledge and Data Engineering, v. 26, n. 9, p. 2250–2267.

[3] H. Wang, M.J. Bah, e M. Hammad, 2019, Progress in Outlier Detection Techniques: A Survey, IEEE Access, v. 7, p. 107964–108000.

# Importance of event detection



## *Event detection initiatives*

---

Anomaly  
detection

Finding unexpected behavior (deviations)

---

Change  
point  
detection

Finding change points

It is related to finding drifts in time series

---

Motif  
detection

Identifying frequent patterns in time series

---

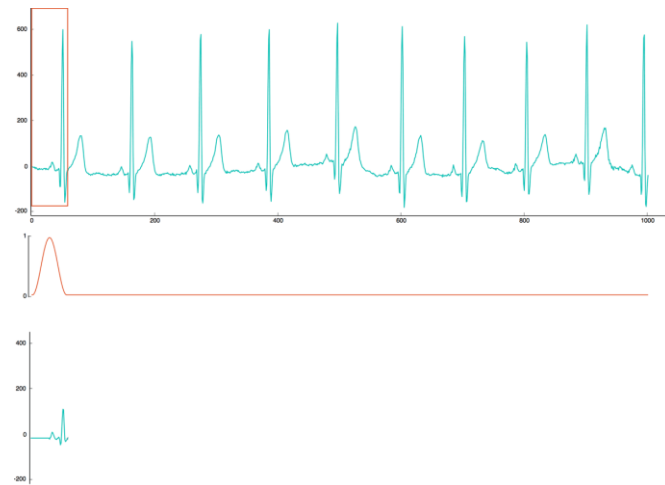
# Anomaly detection (distribution analysis)

- Statistical analysis
  - Differentiation (backshift operator)
  - Residuals from moving average
  - Residuals from filters (Kalman)
- Model adjustment
  - Residuals from decomposed signal
  - Residuals from linear models (regression)
  - Residuals from autoregressive models (ARIMA)
  - Residuals from volatility models (GARCH)
  - Residuals from machine learning models
- Clustering of subsequences
  - Distribution analysis over difference between subsequences and centroids
  - DBScan
- Time series decomposition
  - Trend
  - Seasonal
  - Fourier transform
  - Wavelets
  - IMF - intrinsic mode function
  - Hilbert-Huang transform

# Change point detection

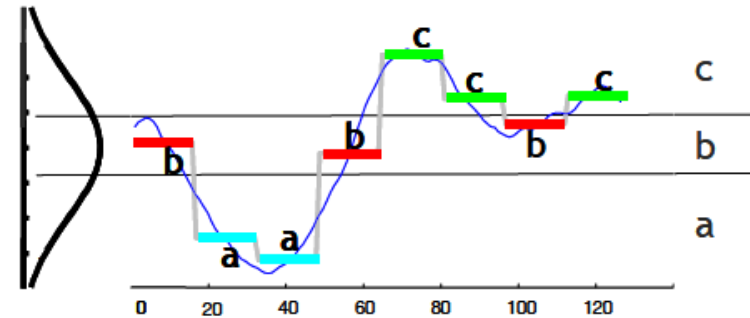
- Seminal change point [1]
- Change Finder [2]

Windowed approach



# Motif discovery

- Indexing
  - Discretization
  - SAX [1]
- Brute force
- Hash-based (random projection) [2]
- Matrix profile [3]



Time Series

X1	X2	X3	X4	X5	X6	X7	X8	X9	X10	X11	X12
----	----	----	----	----	----	----	----	----	-----	-----	-----

X2	X3	X4	X5
----	----	----	----

Distances

D1,2	...	...	...	...	...	...	...	...	...	...	...
------	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----

Matrix Profile Distances

D1	...	...	...	...	...	...	...	...	...	...	...
----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----

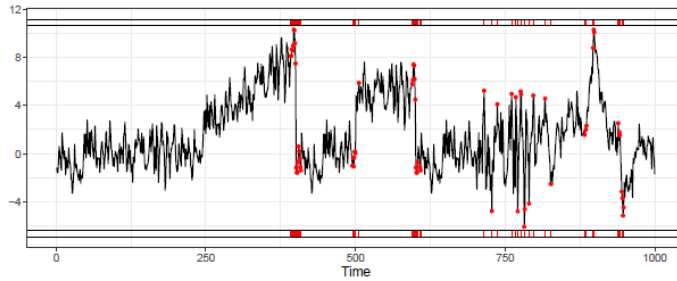
[1] J. Lin, E. Keogh, L. Wei, and S. Lonardi, "Experiencing SAX: a novel symbolic representation of time series," *Data Mining and Knowledge Discovery*, vol. 15, no. 2, pp. 107–144, 2007

[2] A. Mueen, "Time series motif discovery: Dimensions and applications," *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery*, vol. 4, no. 2, pp. 152–159, 2014

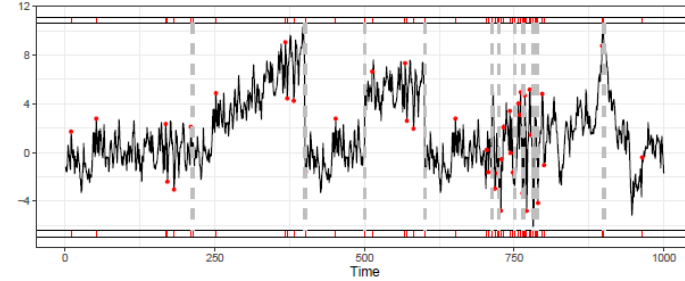
[3] M. Linardi, Y. Zhu, T. Palpanas, and E. Keogh, 2020, Matrix profile goes MAD: variable-length motif and discord discovery in data series, *Data Mining and Knowledge Discovery*, v. 34, n. 4, p. 1022–1071.

(\*) <https://towardsdatascience.com/introduction-to-matrix-profiles-5568f3375d90>

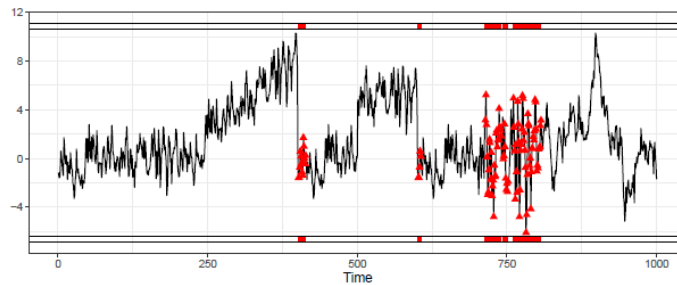
# The many faces of event detection



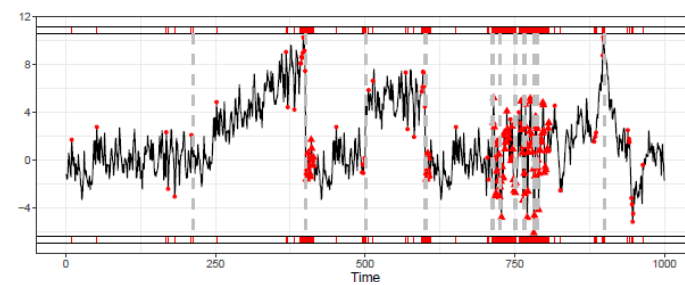
Method A: trend anomalies



Method B: trend anomalies & change points



Method C: volatility anomalies



Methods A,B & C:  
trend anomalies, volatility anomalies and  
change points



# Metrics for event detection

- Classifier Accuracy: percentage of test set tuples that are correctly classified

- $accuracy = \frac{TP+TN}{All}$

- $precision = \frac{TP}{TP+FP}$

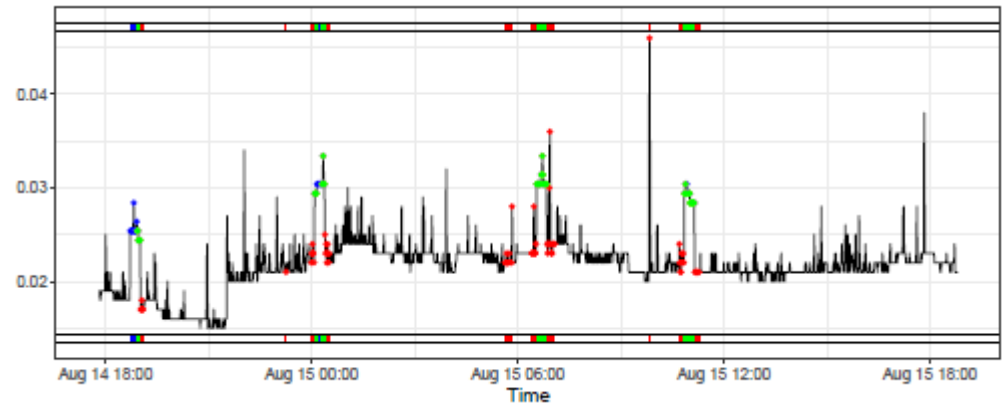
- $recall = \frac{TP}{TP+FN}$

- $F_1 = \frac{2 \cdot precision \cdot recall}{precision + recall}$

- ROC Curve

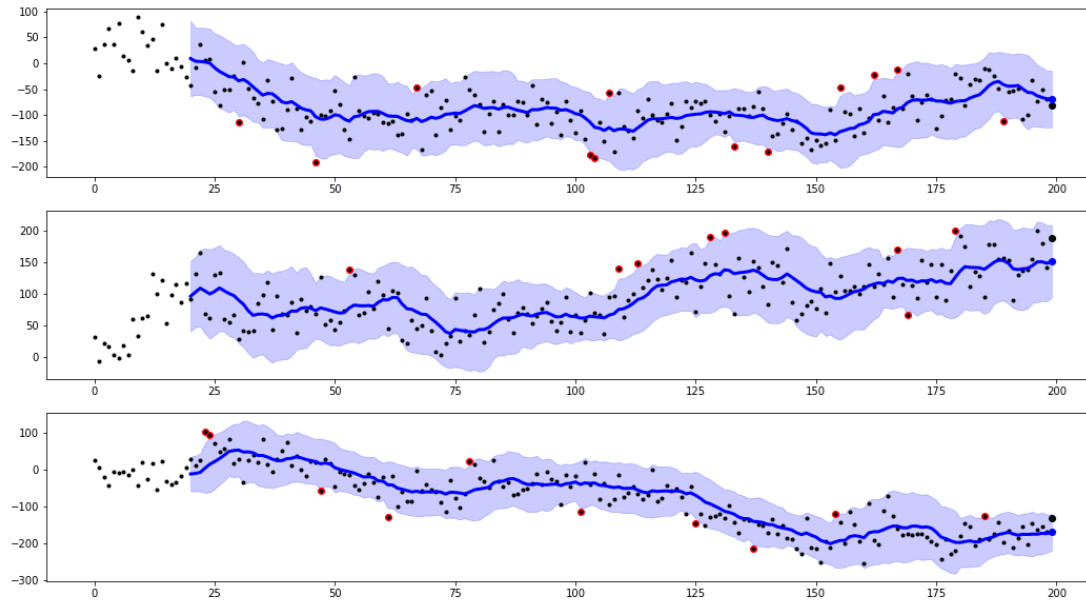
Confusion Matrix (CM)

Predicted Actual	$\hat{E}$	$\neg\hat{E}$
E	TP	FN
$\neg E$	FP	TN

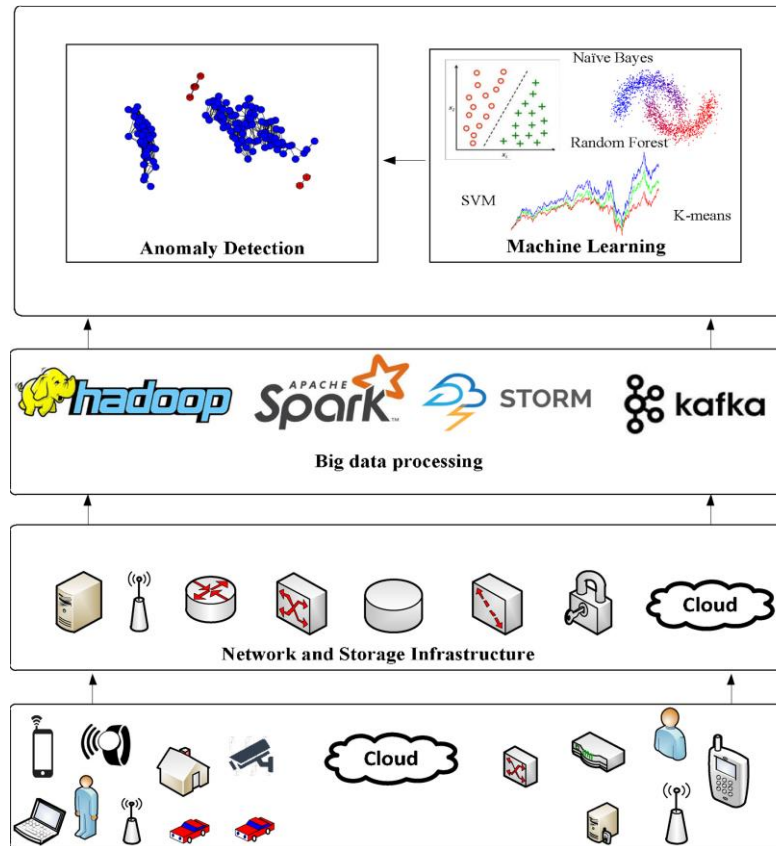


# Online event detection

- Handles streaming time series

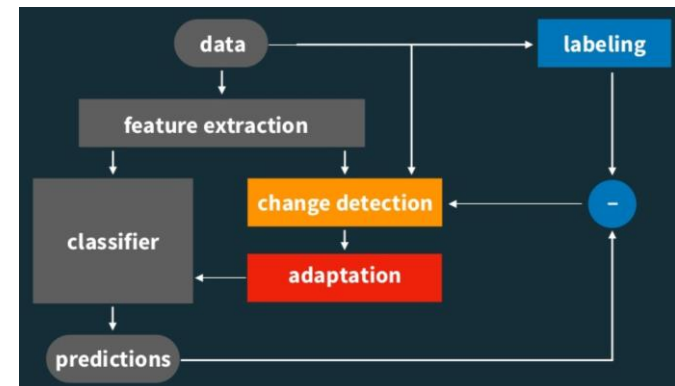
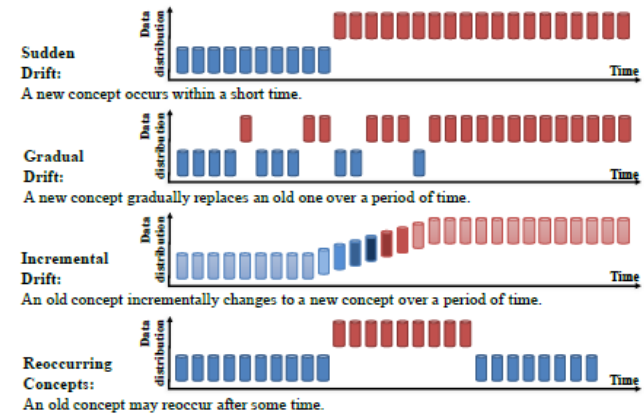
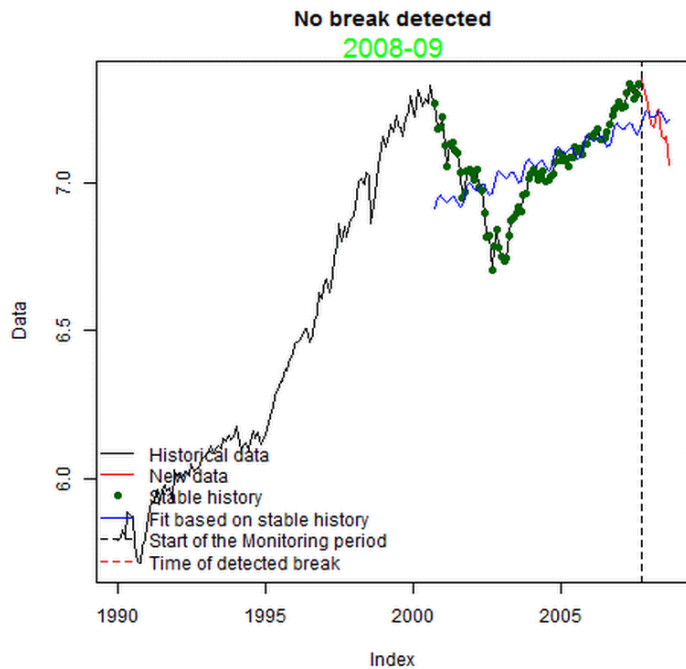


# Online event detection infrastructure



# Online change-point detection

- Detection occurs incrementally

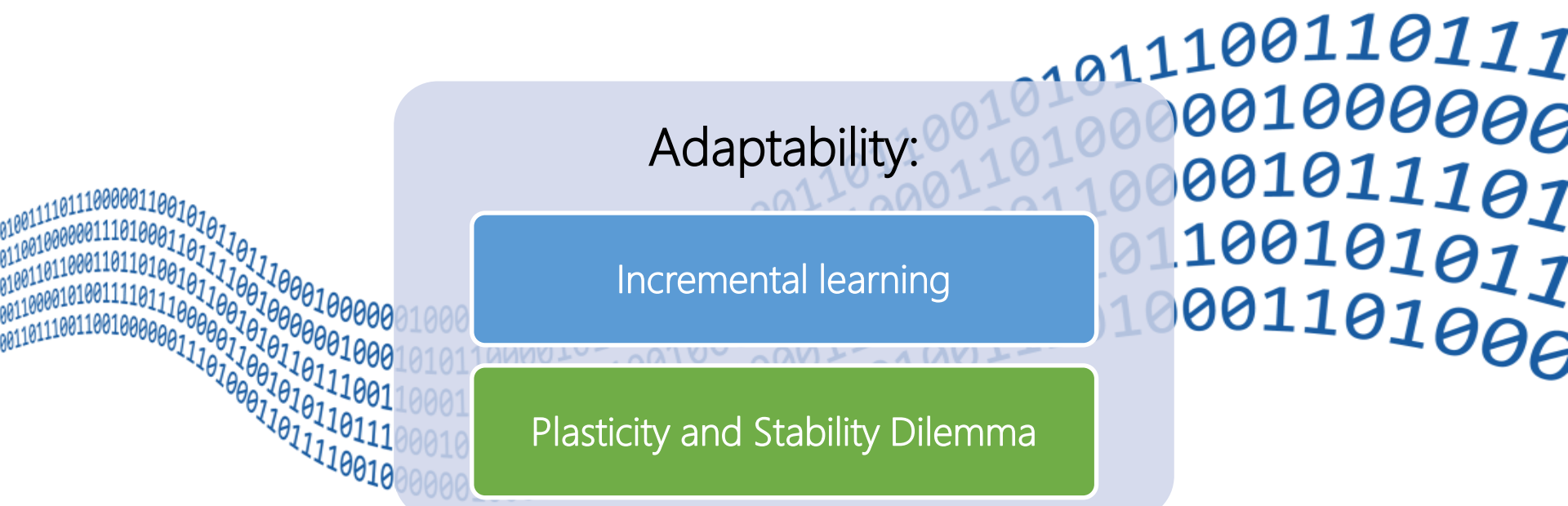


# Online event detection challenges (when to adapt)

Adaptability:

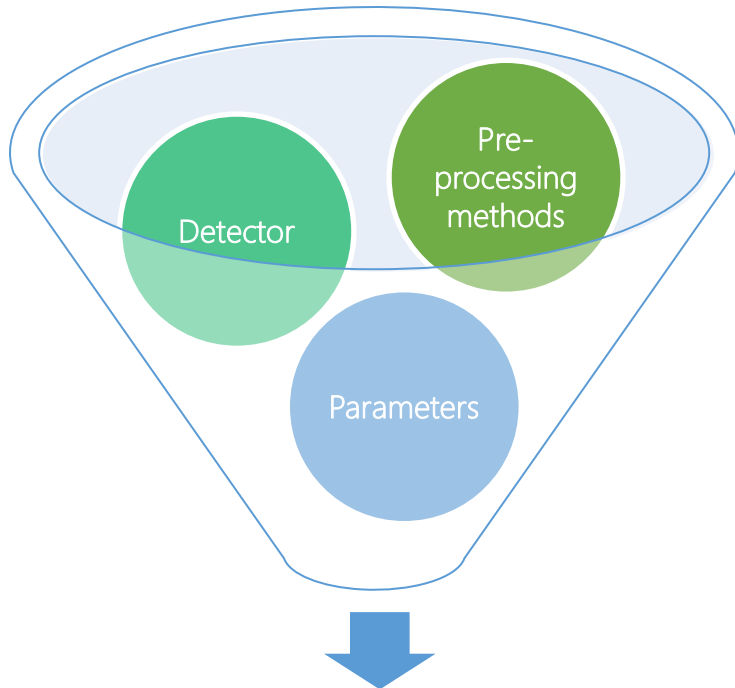
Incremental learning

Plasticity and Stability Dilemma



# Online event detection challenges (too many methods)

Myriad of event detection  
methods (detectors)



Choice of appropriate detectors/parameters for event detection is a challenge

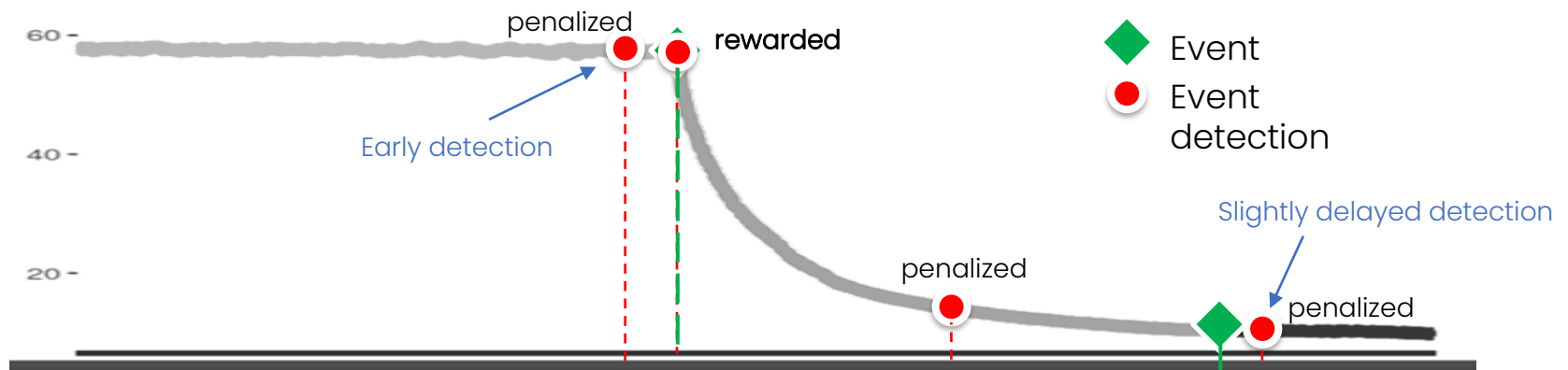
Directly related to initial assumptions about the behavior and statistical properties of data

The nature of the events observed is often unknown

Detectors specialized in a type of event may disregard the occurrence of another, or even misidentify them

# Online event detection challenges (metrics)

- Traditional scoring methods, such as precision and recall, don't suffice for evaluating online event detection performance.
  - They do not incorporate time and do not reward early detection.
  - True positives are rewarded. All other results are "harshly" and equally punished.



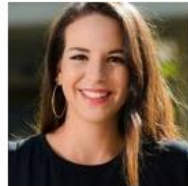
## *What type of disease outbreaks?*

- Doenças Relacionadas ao Saneamento Ambiental Inadequado (DRSAI)
- Internações Hospitalares por Doenças Imunopreveníveis: CID-10 por Doenças Imunopreveníveis
- Taxa de Internação por infecção respiratória aguda de menores de cinco anos de idade (IRA5)



# Data Analytics Lab Team

## Doutorado



**Lais Baroni**  
(CEFET/RJ)



**Leonardo Carvalho**  
(CEFET/RJ)



**Rebecca Salles**  
(CEFET/RJ)

## Mestrado



**Antônio Mello**  
(CEFET/RJ)



**Arthur Garcia**  
(CEFET/RJ)



**Cristiane Gea**  
(CEFET/RJ)



**Diego Salles**  
(CEFET/RJ)



**Flávia Rocha**  
(CEFET/RJ)



**Janio Lima**  
(CEFET/RJ)



**Jéssica de Souza**  
(CEFET/RJ)

# Other researches

- An Analysis of Malaria in the Brazilian Legal Amazon Using Divergent Association Rules
- Estimation of COVID-19 Under-Reporting in the Brazilian States Through SARI
- Neonatal mortality rates in Brazilian municipalities: from 1996 to 2017



Estimation of COVID-19 Under-Reporting in the Brazilian States Through SARI

Rafaela Paula<sup>1</sup>, Luis Baroni<sup>2</sup>, Marcel Pedrosa<sup>3</sup>, Rebeca Salles<sup>4</sup>, Luciano Escobar<sup>5</sup>, Carlos de Souza<sup>6</sup>, Raphael de Freitas Salazar<sup>7</sup>, Jorge Soares<sup>8</sup>, Rafael Coutinho<sup>9</sup>, Fabio Porto<sup>9</sup>, Eduardo Ogasawara<sup>9</sup>

Received: 10 December 2020 / Accepted: 4 March 2021  
© The Author(s) 2021

**Abstract**  
Due to its impact, COVID-19 has been attracting the academy to search for curing, mitigating, or controlling it. It is believed that under-reporting is a relevant factor in determining the actual mortality rate and, if not considered, can cause specific case misinterpretation. Therefore, this work aims to estimate the under-reporting of cases and deaths of COVID-19 in Brazilian states using data from the Instituto Brasileiro de Geografia e Estatística (IBGE) target notification of Severe Acute Respiratory Infection (SARI). The methodology is based on the combination of data analysis (event detection, metrics and time series modeling (metrics and novelty concepts) over hospitalized SARI cases). The estimate of real cases of the disease, called mortality, is calculated by comparing the difference in SARI cases in 2020 (after COVID-19) with the total expected cases in recent years (2016–2019). The expected cases are derived from a seasonal exponential moving average. The results show that under-reporting rates vary significantly between states and that there are no general patterns for states in the same region in Brazil. The states of Minas Gerais and Mato Grosso have the highest rates of under-reporting of cases. The rate of under-reporting of deaths is high in the Rio Grande do Sul and the Mato Grosso. This work can be highlighted for the combination of data analysis and time series modeling. Our calculation of under-reporting rates based on SARI is conservative and better characterized by health than for cases.

**Keywords** COVID-19 · Under-reporting · SARI · Time series modeling · Data analysis

<sup>1</sup> Luis Baroni  
lbaroni@ufpr.br

<sup>2</sup> Federal Center for Technological Education of Rio de Janeiro, CEFET/RJ, Rio de Janeiro, Brazil

<sup>3</sup> Fundação de Amparo à Pesquisa do Estado de São Paulo, FAPESP, São Carlos, Brazil

<sup>4</sup> Instituto de Física de São Carlos, UFSCAR, São Carlos, Brazil

<sup>5</sup> Instituto de Física de São Carlos, UFSCAR, São Carlos, Brazil

<sup>6</sup> Instituto de Física de São Carlos, UFSCAR, São Carlos, Brazil

<sup>7</sup> Instituto de Física de São Carlos, UFSCAR, São Carlos, Brazil

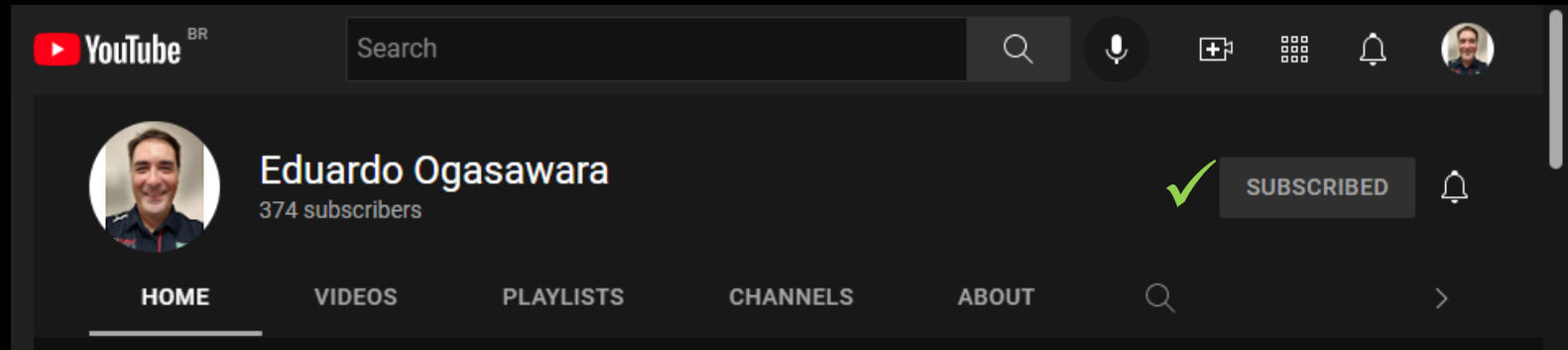
<sup>8</sup> Instituto de Física de São Carlos, UFSCAR, São Carlos, Brazil

<sup>9</sup> Instituto de Física de São Carlos, UFSCAR, São Carlos, Brazil

Published online: 14 March 2021

# Novidades

Inscreva-se em: <https://eic.cefet-rj.br/~eogasawara/youtube>



The image shows a screenshot of a YouTube channel page. At the top left is the YouTube logo with 'BR' next to it. To its right is a search bar with the text 'Search'. Further right are icons for search, voice search, live streaming, a grid of icons, a notification bell, and a profile picture. Below this is the channel header for 'Eduardo Ogasawara', which includes a circular profile picture, the name 'Eduardo Ogasawara', and '374 subscribers'. To the right of the header is a green checkmark, a 'SUBSCRIBED' button, and a notification bell. Below the header is a navigation bar with the following options: 'HOME' (underlined), 'VIDEOS', 'PLAYLISTS', 'CHANNELS', 'ABOUT', a search icon, and a right-pointing arrow.

