

Practical Machine Learning for Streaming Data

ACM SIGKDD Tutorial (hands-on)
2024

Heitor Murilo Gomes¹, Albert Bifet^{2,3}

[https://adaptive-machine-learning.github.io/
kdd2024_ml_for_streams/](https://adaptive-machine-learning.github.io/kdd2024_ml_for_streams/)



[1] Victoria University of Wellington, New Zealand, [2] University of Waikato, New Zealand,
[3] TELECOM Paris, LCTI, France.

**Partially and delayed labeled
data topics**

Partially and delayed labeled data

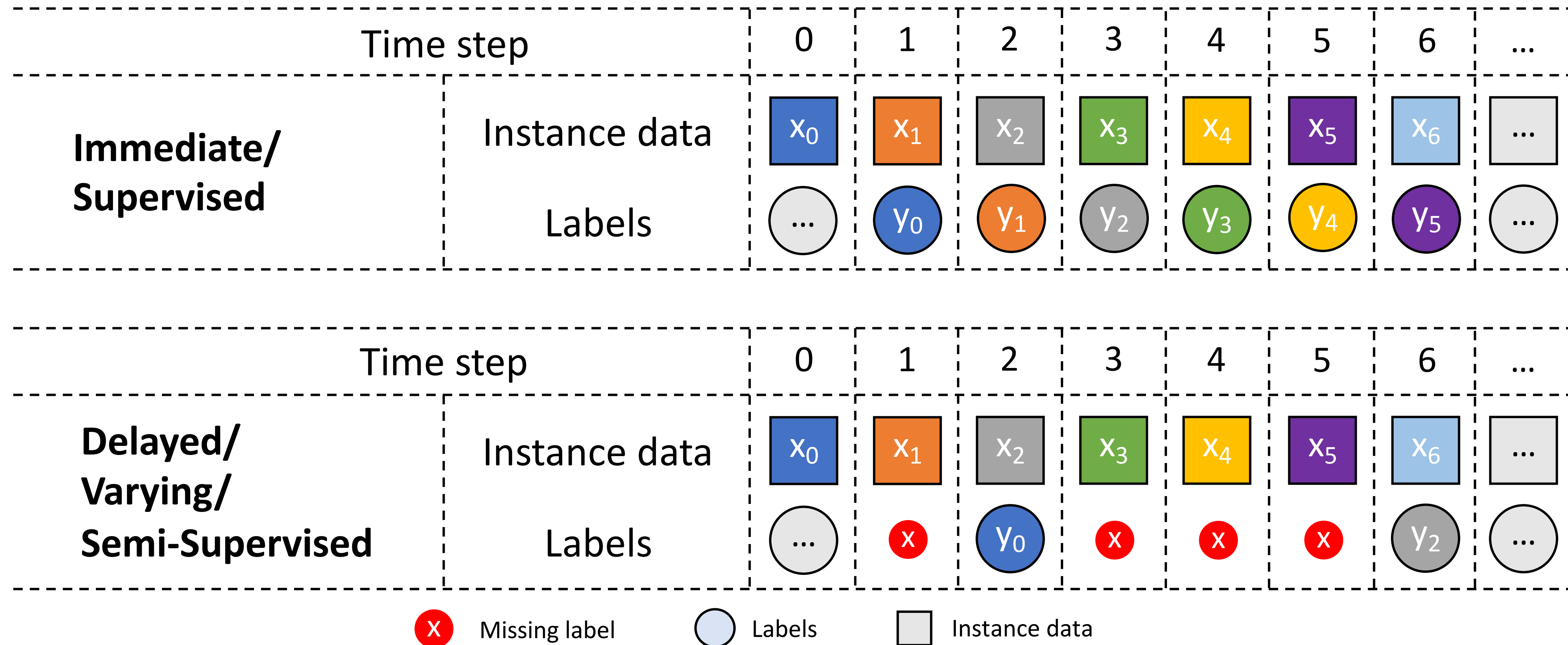
Partially labeled: Not all instances from the stream will be labeled

Delayed labeled: The label may take a while to arrive

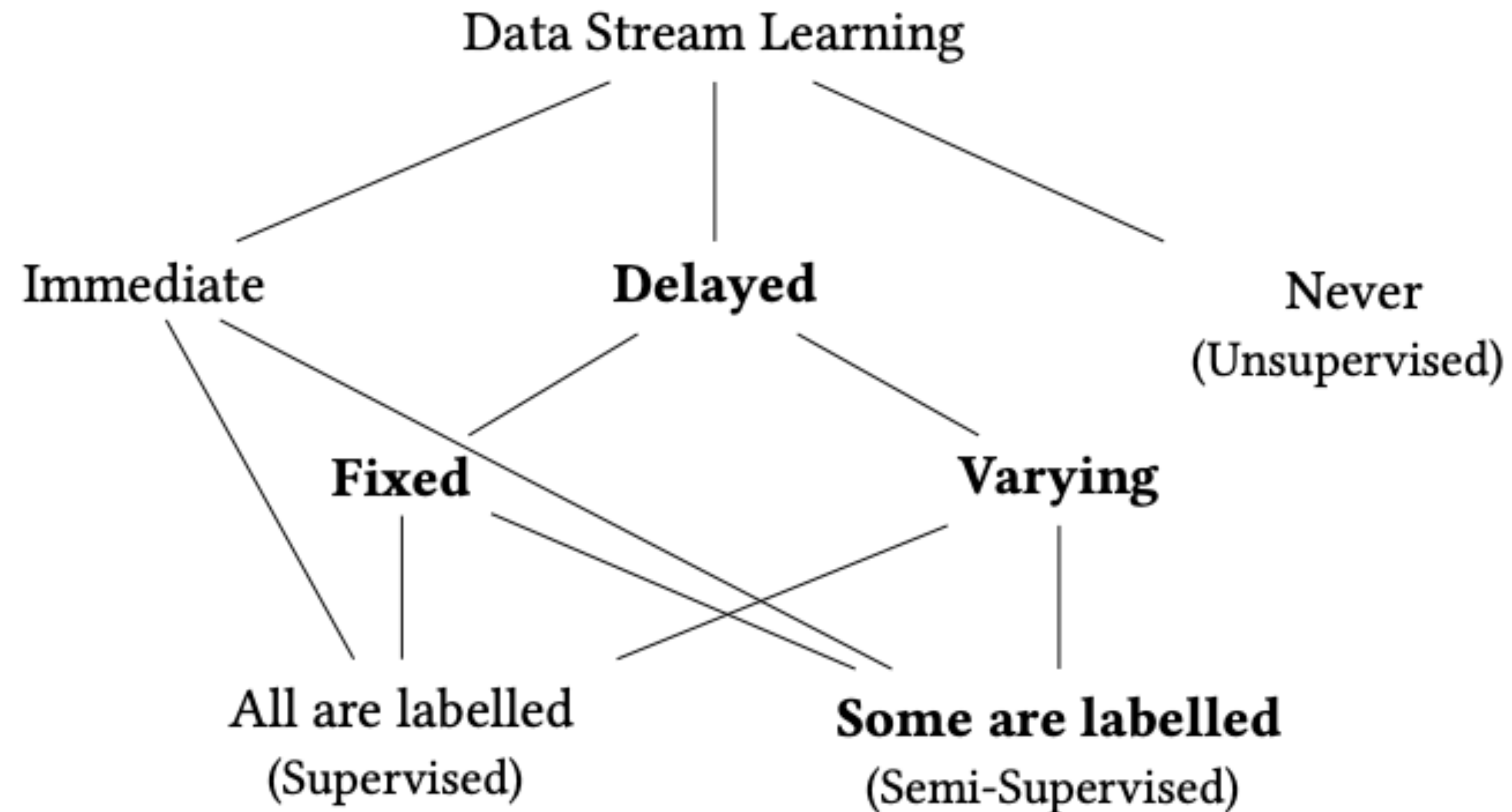
Immediately labeled*: The class label arrives before the next input data to be predicted

Data streams are often associated with large amounts of data. However, most of this data is unlabelled

Partially and delayed labeled data



Learning according to labels arrival time



Semi-supervised learning

Supervised Learning + Unsupervised Learning

*“Typically, **semi-supervised learning** algorithms attempt to **improve performance in one of these two tasks by utilizing information generally associated with the other**. For instance, when tackling a classification problem, additional data points for which the label is unknown might be used to aid in the classification process. For clustering methods, on the other hand, the learning procedure might benefit from the knowledge that certain data points belong to the same class.” [Engelen and Hoos, 2020]*

Related Problems (SSL in DS)

- Active learning
- Transductive learning
- Weakly multi-labelled data
- Initially Labeled Streaming Environment
- Few-shot learning
- Concept evolution
- Online Continual learning

* See section 2.1 from reference

SSL adapted to Streams

- A practical approach to **adapt batch** techniques to streaming is through **sliding or tumbling windows**
- A general strategy for SSL is to use unlabelled examples to **build a representation of the input data**, and then use this representation as input to a model for obtaining predictions.
- The idea is that **an improved representation will lead to improved predictions**; and since representation learning can naturally be an unsupervised task, training labels are not required
- Examples: Restricted Boltzmann machines (RBMs), Auto-encoders, Cluster representations

Evaluation of partially labeled data

- **Labeling ratio**, 1%, 5%, ..., if too much data is labeled, then it is likely that a supervised algorithm would prevail
- **Delay simulation**, often simulated based on number of instances
- **Comparison against SSL and supervised strategies**

Clustering

Clustering Data Streams

Clustering data streams refers to grouping instances into clusters as the data continuously flows in, which normally includes **two phases**:

1. Online Step

- a. Micro-Cluster*** **Formation**: Incoming data points are incrementally processed and assigned to micro-clusters.
- b. Micro-Cluster Maintenance**: The micro-clusters are periodically updated as new data arrives. This includes adjusting the micro-cluster centroids and merging or splitting clusters based on defined thresholds.

2. Offline Step

Periodically or upon request, micro-clusters are aggregated into macro-clusters (or simply clusters) to provide a higher-level view of the data.

* **Micro-clusters** are small, temporary “clusters” that capture local density information and are typically represented by statistical summaries like centroid, weight, and radius.

CluStream

1. Online Step

For each new **instance i** that arrives:

i is **absorbed** by a micro-cluster

i **starts** a new **micro-cluster** of its own

Delete oldest micro-cluster

Merge two of the **oldest** micro-cluster

2. Offline Step

Apply k-means using micro-clusters as instances

Anomaly Detection

Anomaly Detection for Data Streams

- Identification of anomalous data in a continuous flow of data
- **Challenges**
 - Adapting to concept drifts without missing out on anomalies
 - Detecting rare anomalies amidst high-volume data streams
 - And more...

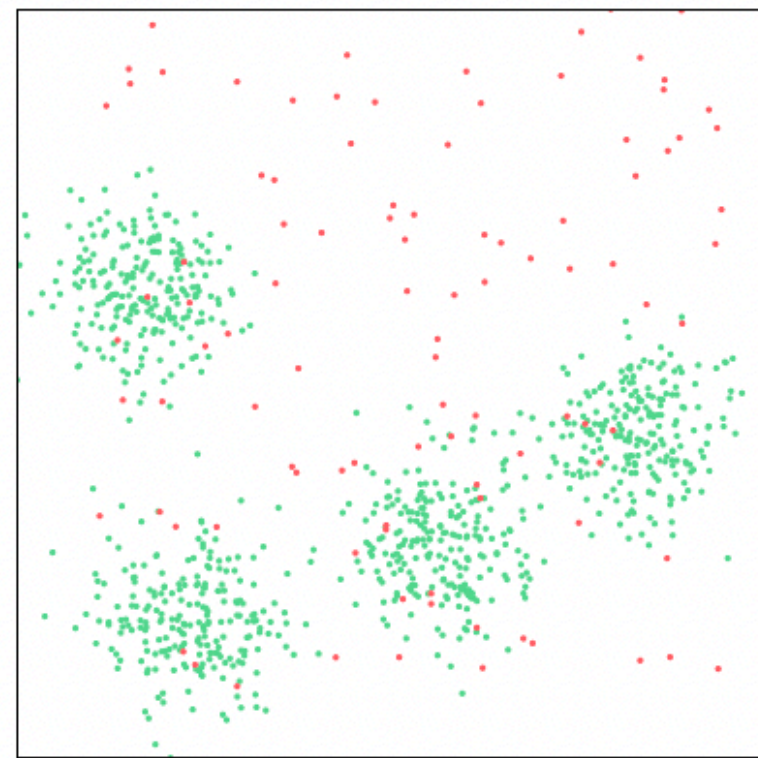
Online Isolation Forest (OIF)

- Inspired by the classic Isolation Forest [1]
- OIF uses a **group of histograms at different levels of detail to capture** the data patterns, with a flexible system that can **learn from new data and gradually forget older data.**

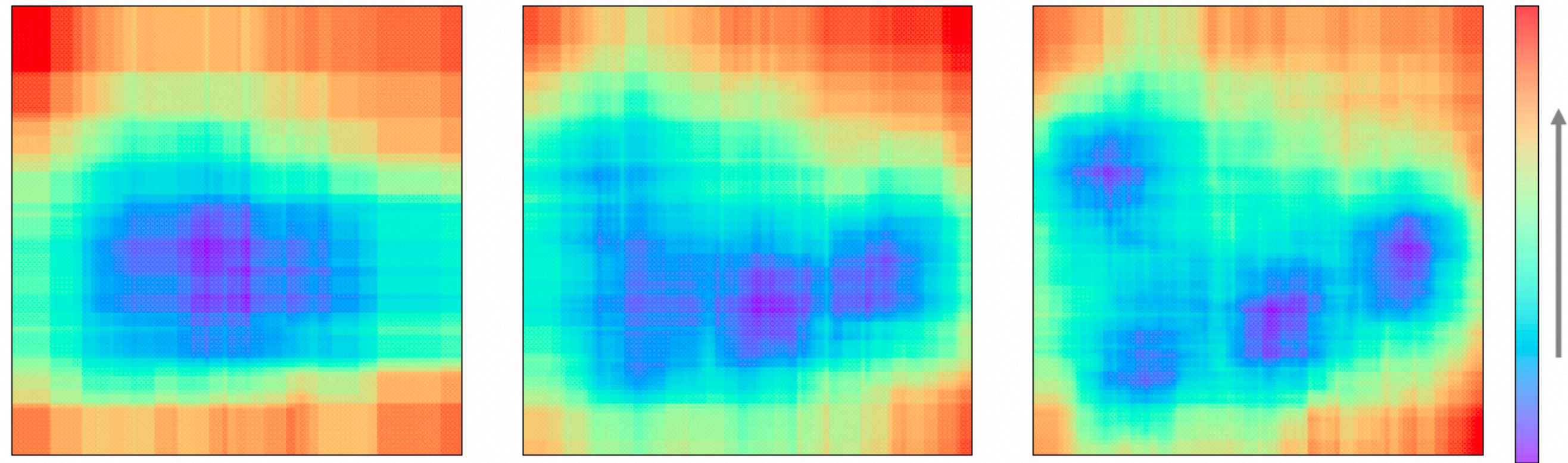
[1] Fei Tony Liu, Kai Ming Ting, and Zhi-Hua Zhou. "Isolation forest." ICDM, 2008.

[2] Filippo Leveni, G W Cassales, B Pfahringer, A Bifet, and G Boracchi. "Online Isolation Forest." ICML, 2024

Online Isolation Forest



(a) Data stream $\mathbf{x}_1, \dots, \mathbf{x}_t \in \mathbb{R}^d$.



(b) Anomaly scores s at different time instants t , from left to right.

- Genuine data (green) are more densely distributed than anomalous data (red)
- OIF processes each data point individually online, assigning an anomaly score to each
- As more data is available, OIF continuously updates and refines the anomaly scores based on the evolving data distribution.

Hands-on example

KDD_2024_advanced.ipynb

and

KDD_2024_solutions.ipynb

Coming up next in 2024

- **Upcoming Tutorials**
 - PAKDD: May 2024 (Taipei, Taiwan) **[done!]**
 - IJCAI: August 2024 (Jeju, South Korea) **[done!]**
 - Kiwi Pycon: August 2024 (Wellington, NZ) **[done!]**
 - KDD: August 2024 (Barcelona, Spain) **[this one]**
 - ECML: September 2024 (Vilnius, Lithuania)
- **CapyMOA next release** (September 2024)



Conclusion

- Streaming data is everywhere — often in a delayed and partially labeled
- ML algorithms for data streams should be **accurate**, **adaptive** and **efficient**
- **CapyMOA** can be easily extended for many stream tasks

Contact: heitor.gomes@vuw.ac.nz



<https://discord.gg/RekJArWKNZ>

Thank you!



<https://github.com/adaptive-machine-learning/CapyMOA>