

# StreamRHF in CappyMOA

Giovanni BENEDETTI DA ROSA  
Cristian Alejandro CHÁVEZ BECERRA

Master 2 Data Science  
Institut Polytechnique de Paris-École Polytechnique  
Data Streaming

# Introduction

# Introduction

- **Data streaming** applications: fraud detection, network monitoring, and real-time recommendation system
- **Anomaly Detections** algorithms identifying interesting or rare events that deviate from normality



Image extracted from [1]

# Problematic

- Data streaming algorithms: memory requirements and faster processing times.
- Anomaly Detections algorithms generally not adapt to the requirements.
- **Implement StreamRHF[2] in the CopyMOA[1] library.**

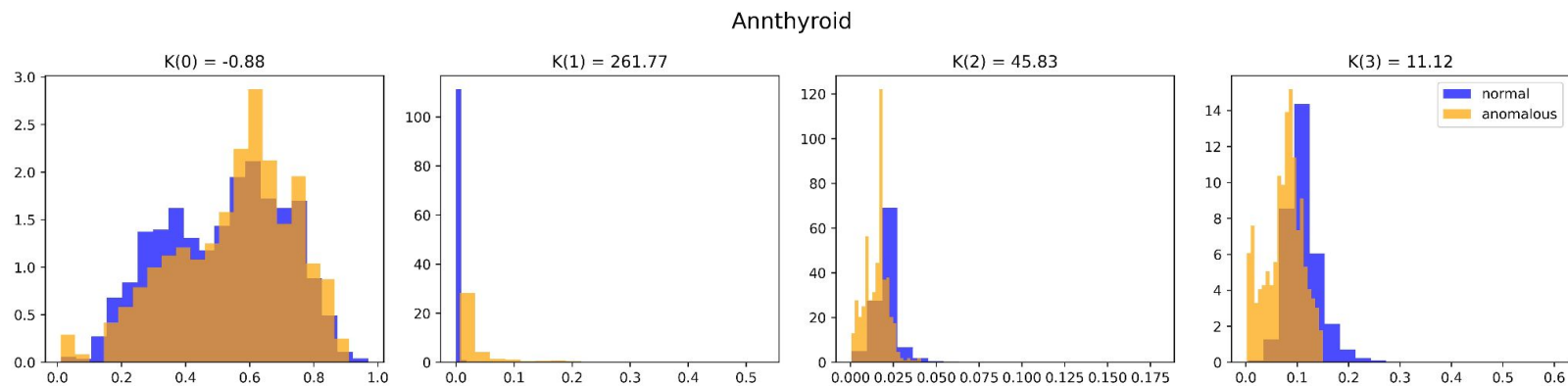


Image extracted from [1]

# Methods

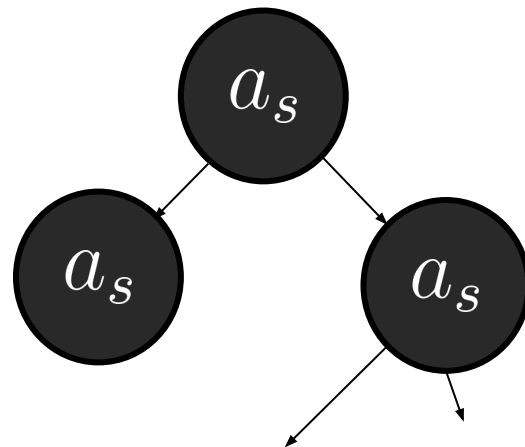
# Methods: RHF

Moors et al[4] described kurtosis as a measure of dispersion and explained that high kurtosis values result from either: (i) occasional values far from the mean in a distribution with most of its probability mass concentrated around the mean, or (ii) significant probability mass located in the distribution's tails.



# Methods: RHF

- Unsupervised Anomaly Detection[3] method
- Data is split across different bins by means: each one of the  $t$  trees with height  $h$  is built by recursively splitting the whole dataset
- Each splitting decision is done selecting an attribute  $a$  to split according to its kurtosis score and a split value randomly



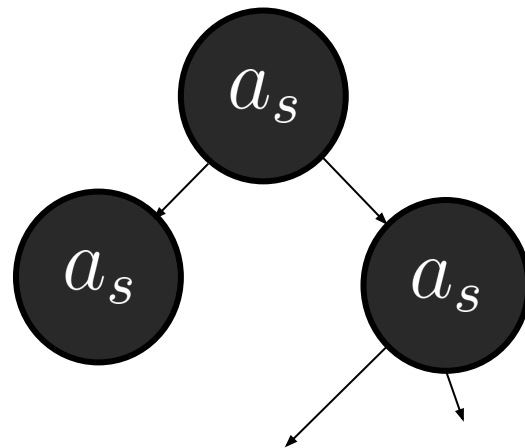
# Methods: RHF

- As defined by Putina et al. Given a dataset  $X \in \mathbb{R}^{n \times d}$  we can define kurtosis split as:

$$1) K_s = \sum_{a=0}^d \log(K(X_a) + 1)$$

$$2) r \sim \text{Uniform}([0, K_s))$$

$$3) a_s = \arg \min \left( k \in [0, d] \left| \sum_{a=0}^k \log(K(X_a) + 1) > r \right. \right)$$





# Methods: RHF

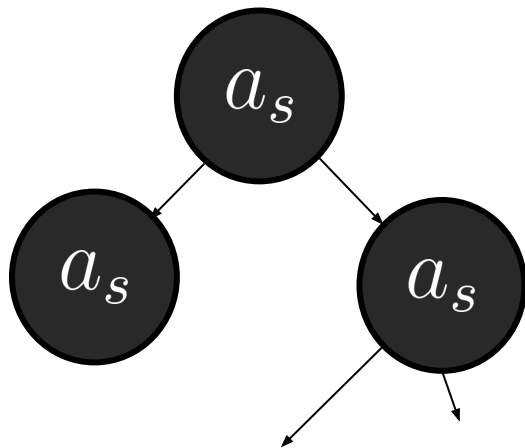
- **Anomaly Score:** Given that  $S_l$  represents the number of instances in each leaf, each instance  $i$  is stored in a tree  $T$ , it's possible to write the information content:

$$w_T(i) = \log \left( \frac{1}{P_{l_i}} \right)$$

where  $P_l$  is the probability associated with leaf  $l$ .

The **anomaly score** is defined as:

$$w_i = \sum_T w_T(I).$$



# Methods: StreamRHF

In streaming of datas, distribution of data is not time invariant, leading to concept drifts.

Nesic et. al[2] define **StreamRHF** uses a window(n) approach:

- 1) Wait for n initial points of the data, than build a tree using RHF and compute the score.
- 2) At each time a new instance arrive, check if the splitting attribute of each node will change; if it's the case: rebuild the subtree; Compute new score.
- 3) When new n data points arrive rebuild RHF, using these data.

# Methods: StreamRHF

- **Time and Space Complexity:** Let  $w$  be the size of the sliding window,  $d$  the number of features,  $h$  the maximum height of every tree, and  $t$  be the total number of trees in our forest[2].
- For the space complexity:  $O(wdt)$
- We rebuild only relatively small subtrees  $\bar{h} < h$  when doing an insertion with height  $\bar{h}$ , which gives a running time of  $O(wdt\bar{h})$  for an update

$$O(wdt\bar{h})$$

# Experimental Evaluation

# Parameters settings

- To compare our implementation, we used the default values of the built-in implemented methods in CappyMOA, using the following parameters:
  1. **HalfSpaceTrees**[5]: window size = 100, number of trees = 25, max depth = 15, anomaly threshold = 0.5.
  2. **Online Isolation Forest**[6]: num trees = 32, max leaf samples = 32, growth criterion = 'adaptive', subsample = 1.0, window size = 2048, branching factor = 2, split = 'axisparallel'.
  3. **Autoencoder** [7]: hidden layer = 2, learning rate = 0.5, threshold = 0.6.
  4. **StreamRHF**: number\_of\_trees = 100, max-height = 5, window\_size: 1% of the dataset.

# Metrics

- Imbalanced classes - intrinsic characteristic of anomaly detections.
- Streaming data the proportion of these anomalies is not constant. While is well known that AUCPR takes to account these for fix outlier proportion than in ROC AUC. Some others work shows the opposite for anomaly detections[8]. We then choose:

- **Average Precision**  $AP = \sum (R_n - R_{n-1})P_n,$

- **AUC** 
$$AUC = \sum_{i=1}^N (FPR_i - FPR_{i-1}) \cdot TPR_i,$$

# Implementation in CapyMOA

- **Three classes implemented:** Node, RandomHistogramForest, StreamRHF
- **Normalization of the scores:** We defined a maximum and minimum value for the scores calculated by the original paper of StreamRHF. Initially, the score is given by the following equation:

$$w_T(i) = \log \left( \frac{1}{P_{l_i}} \right) = \log \left( \frac{n}{|S(l)|} \right)$$

## Implementation in CapyMOA: Normalization of the scores

$$\text{min\_score} = \log \left( \frac{n}{n} \right) = \log(1) = 0$$

$$\text{max\_score} = \log \left( \frac{n}{1} \right) = \log(n)$$

$$\text{normalized\_score} = \frac{w_T(i) - \text{min\_score}}{\text{max\_score} - \text{min\_score}}$$



# Results: Comparison using the AP score

Dataset	StreamRHF	Autoencoder	HalfSpaceTrees	OnlineIsolationForest
abalone	0.237 $\pm$ 0.024	0.116	<b>0.531</b>	0.033
annthyroid	0.396 $\pm$ 0.019	0.086	<b>0.577</b>	0.079
kdd ftp	0.280 $\pm$ 0.008	0.157	<b>0.439</b>	0.196
magicgamma	<b>0.623 <math>\pm</math> 0.008</b>	0.352	0.272	0.287
mammography	<b>0.219 <math>\pm</math> 0.021</b>	0.157	0.184	0.049
mnist	<b>0.303 <math>\pm</math> 0.006</b>	0.092	0.117	0.085
musk	<b>0.771 <math>\pm</math> 0.061</b>	0.032	0.021	0.495
satellite	<b>0.593 <math>\pm</math> 0.012</b>	0.015	0.317	0.518
satimages	<b>0.856 <math>\pm</math> 0.020</b>	0.012	0.256	0.012
shuttle_odds	<b>0.858 <math>\pm</math> 0.044</b>	0.072	0.548	0.040
spambase	0.478 $\pm$ 0.038	0.507	0.282	<b>0.568</b>
thyroid	<b>0.544 <math>\pm</math> 0.031</b>	0.103	0.559	0.014
Mean	<b>0.513 <math>\pm</math> 0.024</b>	0.142	0.342	0.198
Median	<b>0.511</b>	0.098	0.300	0.082

# Results: Comparison using the AP score

1. We achieved very similar results compared to the ones shown in the original paper in the datasets we used using this score (the score used in the original paper)

$$AP = \sum_n (R_n - R_{n-1}) P_n,$$

2. StreamRHF generally outperforms the default configurations of the built-in CapyMOA models.
3. In the case of very small datasets, such as the abalone dataset, HalfSpaceTrees surpasses StreamRHF.

# Results: Comparison using the AUC score

Dataset	StreamRHF	Autoencoder	HalfSpaceTrees	OnlineIsolationForest
abalone	0.862 $\pm$ 0.005	0.855	<b>0.963</b>	0.318
annthyroid	0.865 $\pm$ 0.006	0.534	<b>0.974</b>	0.516
kdd ftp	0.423 $\pm$ 0.022	0.124	<b>0.662</b>	0.350
magicgamma	<b>0.701 <math>\pm</math> 0.010</b>	0.500	0.382	0.396
mammography	0.858 $\pm$ 0.004	0.766	<b>0.912</b>	0.749
mnist	<b>0.809 <math>\pm</math> 0.005</b>	0.500	0.558	0.499
musk	<b>0.980 <math>\pm</math> 0.007</b>	0.500	0.026	0.612
satellite	0.939 $\pm$ 0.009	0.500	<b>0.975</b>	0.541
satimages	<b>0.989 <math>\pm</math> 0.003</b>	0.500	0.965	0.494
shuttle_odds	<b>0.994 <math>\pm</math> 0.003</b>	0.500	0.899	0.101
spambase	0.636 $\pm$ 0.034	<b>0.678</b>	0.271	0.496
thyroid	0.977 $\pm$ 0.001	0.823	<b>0.979</b>	0.168
Mean	<b>0.839 <math>\pm</math> 0.011</b>	0.556	0.725	0.423
Median	<b>0.862</b>	0.500	0.662	0.396

# Results: Comparison using the AUC score

1. We decided to use the AUC metric since is a standard metric for Anomaly Detection models that was not used in the original paper.

$$AUC = \sum_{i=1}^N (FPR_i - FPR_{i-1}) \cdot TPR_i,$$

2. StreamRHF doesn't always outperform other CapyMOA models in terms of AUC, but it delivers competitive results, often close to the best model, HalfSpaceTrees
3. Its average and median AUC values are still higher, showing its overall reliability.

# Results: Insertion time

Dataset	StreamRHF	Autoencoder	HalfSpaceTrees	OnlineIsolationForest
abalone	117.338 $\pm$ 1.076	0.310 $\pm$ 0.015	<b>0.098 <math>\pm</math> 0.026</b>	5.642 $\pm$ 0.044
annthyroid	122.116 $\pm$ 1.671	0.303 $\pm$ 0.001	<b>0.032 <math>\pm</math> 0.005</b>	10.633 $\pm$ 0.023
kdd ftp	112.115 $\pm$ 1.050	0.298 $\pm$ 0.000	<b>0.032 <math>\pm</math> 0.001</b>	10.603 $\pm$ 0.020
magicgamma	202.003 $\pm$ 7.166	0.305 $\pm$ 0.001	<b>0.033 <math>\pm</math> 0.010</b>	13.816 $\pm$ 0.032
mammography	134.445 $\pm$ 2.414	0.304 $\pm$ 0.003	<b>0.031 <math>\pm</math> 0.003</b>	13.813 $\pm$ 0.018
mnist	398.449 $\pm$ 5.247	0.370 $\pm$ 0.001	<b>0.032 <math>\pm</math> 0.000</b>	13.300 $\pm$ 0.023
musk	547.362 $\pm$ 12.781	0.422 $\pm$ 0.002	<b>0.045 <math>\pm</math> 0.024</b>	8.237 $\pm$ 0.040
satellite	206.419 $\pm$ 4.093	0.324 $\pm$ 0.000	<b>0.032 <math>\pm</math> 0.000</b>	9.457 $\pm$ 0.026
satimages	205.343 $\pm$ 6.415	0.323 $\pm$ 0.001	<b>0.032 <math>\pm</math> 0.000</b>	11.490 $\pm$ 0.033
shuttle.odds	189.577 $\pm$ 6.058	0.304 $\pm$ 0.002	<b>0.029 <math>\pm</math> 0.000</b>	12.556 $\pm$ 0.275
spambase	177.934 $\pm$ 2.605	0.338 $\pm$ 0.001	<b>0.032 <math>\pm</math> 0.000</b>	7.627 $\pm$ 0.031
thyroid	112.401 $\pm$ 0.947	0.303 $\pm$ 0.001	<b>0.037 <math>\pm</math> 0.007</b>	8.647 $\pm$ 0.021

Table 3: Insertion times (in ms) for different datasets and models, with the smallest values highlighted in bold.

## Results: Insertion time

1. Taking into account the importance of time in a data streaming framework we decided to compare the insertion time of the models.
2. StreamRHF takes significantly longer than other models, which can be attributed to several factors: default number of trees, maximum height of the trees and language implementation (cython in the original paper)
3. We may reduce the insertion time by reducing the number of trees and/or maximum height of the trees but it is a trade-off in performance and speed

## Results: CapyMOA check

Dataset	Original Paper Procedure (AP)	CapyMOA Procedure (AP)
abalone	$0.233 \pm 0.035$	$0.123 \pm 0.020$
annthyroid	$0.401 \pm 0.035$	$0.354 \pm 0.028$
kdd ftp	$0.280 \pm 0.011$	$0.169 \pm 0.036$
thyroid	$0.548 \pm 0.044$	$0.444 \pm 0.017$
Mean	$0.365 \pm 0.031$	$0.272 \pm 0.012$
Median	0.341	0.259

## Results: CapyMOA check

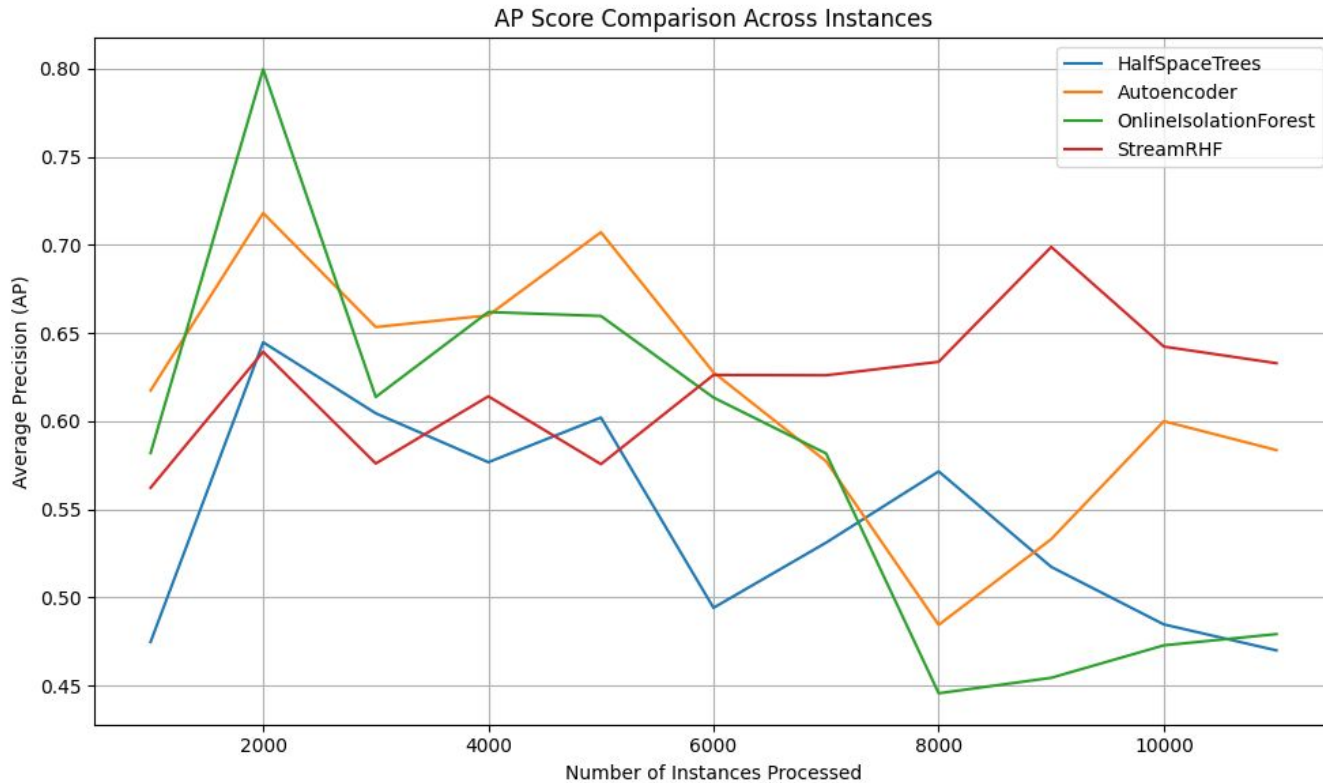
1. In the original paper, the procedure to evaluate the model the authors did for every instance was to first insert the instance and then score the instance. The other way around according to CapyMOA tutorials.

$$w_T(i) = \log \left( \frac{1}{P_{l_i}} \right) = \log \left( \frac{n}{|S(l)|} \right)$$

2. This procedure reduces the anomaly score for the all the instances but in different proportions.
3. As can be seen in table, the original paper procedure leads to better results for the AP score.



# Demo



# Conclusion

# Conclusion

- StreamRHF is a novel method that demonstrates promising results and competes with state-of-the-art data streaming models.
- We were able to reproduce results similar to those in the original paper.
- We successfully integrated our implementation with CappyMOA.
- Additionally, we established maximum and minimum scores to facilitate its integration into CappyMOA.
- Our implementation in CappyMOA provides results comparable to the original implementation in the .py file.
- However, the main drawback of our implementation is the computational time, which is impacted by the high number of trees and the use of a less efficient programming language.

# Code Availability

Our GitHub repository for the project can be found at:

<https://github.com/AlejandroUN/Stream-Random-Histogram-Forest>

The from-scratch implementation Python file of the algorithm StreamRHF can be found in the following link in our GitHub repository:

[https://github.com/AlejandroUN/Stream-Random-Histogram-Forest/blob/main/src/capymoa/anomaly/\\_stream\\_rhf.py](https://github.com/AlejandroUN/Stream-Random-Histogram-Forest/blob/main/src/capymoa/anomaly/_stream_rhf.py)

The demo showing the implementation properly integrated with CapyMOA using CapyMOA models and plots to compare can be found in the following notebook:

[https://github.com/AlejandroUN/Stream-Random-Histogram-Forest/blob/main/src/capymoa/anomaly/streamrhf/notebook\\_presentation.ipynb](https://github.com/AlejandroUN/Stream-Random-Histogram-Forest/blob/main/src/capymoa/anomaly/streamrhf/notebook_presentation.ipynb)

# References

1. Gomes, H. M., Bifet, A. (2024). Practical Machine Learning for Streaming Data. ACM Digital Library, 6418-6419. Available: <https://doi.org/10.1145/3637528.3671442>
2. Stefan Nesic, Andrian Putina, Maroua Bahri, Alexis Huet, Jose Manuel, et al.. STREAMRHF: Tree-Based Unsupervised Anomaly Detection for Data Streams. AICCSA 2022 - 19th ACS/IEEE International Conference on Computer Systems and Applications, Dec 2022, Abu Dhabi, United Arab Emirates. fffhal-03948938f
3. A. Putina, M. Sozio, D. Rossi, and J. M. Navarro, "Random histogram forest for unsupervised anomaly detection," in International Conference on Data Mining (ICDM). IEEE, 2020, pp. 1226–1231.
4. J. J. A. Moors, "The meaning of kurtosis: Darlington reexamined", The American Statistician, vol. 40, no. 4, pp. 283–284, 1986.
5. Swee Chuan Tan, Kai Ming Ting, and Tony Fei Liu. Fast anomaly detection for streaming data. In International Joint Conference on Artificial Intelligence (IJCAI), pages 1469–1495, 2017.
6. Filippo Leveni, Guilherme Weigert Cassales, Bernhard Pfahringer, Albert Bifet, and Giacomo Boracchi. Online Isolation Forest. In International Conference on Machine Learning (ICML), Proceedings of Machine Learning Research (PMLR), 2024.
7. Richard Hugh Moulton, Herna L. Viktor, Nathalie Japkowicz, and Jo˜ao Gama. Contextual one-class classification in data streams. arXiv preprint arXiv:1907.04233, 2019.
8. Minjae Ok, Simon Kluttermann, and Emmanuel Mˆuller, "Exploring the Impact of Outlier Variability on Anomaly Detection Evaluation Metrics," arXiv preprint, arXiv:2409.15986v1 [cs.LG], 24 Sep 2024, License: CC BY 4.0.
9. Bifet, Albert, et al. Machine Learning for Data Streams: with Practical Examples in MOA. The MIT Press, 2018. DOI: <https://doi.org/10.7551/mitpress/10654.001.0001>.