

SDOstream: Low-Density Models for Streaming Outlier Detection

Alexander Hartl, Félix Iglesias and Tanja Zseby *

TU Wien - Institute of Telecommunications
Gußhausstraße 25, 1040 Vienna - Austria

Abstract. Data commonly changes over time. Algorithms for anomaly detection must therefore be adapted to overcome the challenges of evolving data. We present SDOstream, a distance-based outlier detection algorithm for stream data that uses low-density models, therefore operating in linear time and avoiding the limitations of sliding windows and instance-based methods. SDOstream is designed to ensure a good integration in applications, hence the definition of “outlier” is not predetermined, but can be decided by the application based on distances to representative point locations. We describe the algorithm and evaluate algorithm performance with several datasets.

1 Introduction

Dynamic data shows peculiarities that reduce the efficacy of traditional machine learning. The main challenge is often referred to as *concept drift*, addressing the evolution of the structures that underlie data. Such temporal changes progressively downgrade the validity of machine learning unless it is adaptive.

Some alternatives to implement adaptiveness are re-fitting or updating models. However, anomaly (or outlier) detection is traditionally faced using instance-based methods like k-neighborhood (KNN) or LOF [1]. As general purpose options, such methods have not been overcome yet in terms of accuracy [2]. On the other hand, their main drawback is the need to consider all previous data for every new data point, a fact that significantly slows down the analysis. In stream data, this problem is partially alleviated by using a *sliding window* (SW), which can be considered a memory length. Moreover, *outlierness* is strictly defined based on non-intuitive parameters, many times being arbitrarily adjusted.

We introduce SDOstream, a simple and elegant outlier detection algorithm for evolving data streams, which operates in $O(n)$ in the number of processed data points. SDOstream builds a model for the data and, hence, avoids the need to store all points in a SW. SWs define a cut-off time, before which data points are not remembered. In SDOstream, impact of past data points is of exponential shape. Hence, it has the potential to retain shapes for a much longer time period. Therefore, whereas in SWs memory length frequently has to be decided based on algorithm limitations, in SDOstream it is defined only based on the application.

Furthermore, in SDOstream the final definition of outlierness is established based on space distances and representative point locations, providing an intuitive and interpretable decision process.

*This work was partly supported by the project MALware cOmmunication in cRitical Infrastructures (MALORI), funded by the Austrian security research programme KIRAS of the Federal Ministry for Transport, Innovation and Technology (bmvit).

2 Related Work

Most methods based on SWs share a strict definition of outlieriness [3]:

Let S be a set of objects, obj an object of S , k a positive integer, and R a positive real number. Then, obj is a distance-based outlier (or, simply, an outlier) if less than k objects in S lie within distance R from obj .

AbstractC [4], Exact- and Approx-STORM [3] and the COD family [5] adopt this definition, mainly differing in their approaches to reduce processing time, e.g., using sampling (Approx-STORM) or micro-clusters (MCOD).

Conceptually different from the previous methods, AnyOut [6] uses a ClusTree model, relying on hierarchical clustering to summarize the data into condensed, evolving clusters. This strategy might face strong complications when data does not match obvious cluster schemes. Also using models, methods based on Random Forests (e.g., [7, 8]) are very powerful due to their speed, robustness, and feature independence. Several approaches require labeled training data or simply establish outlieriness based on estimations of the data mass. xStream [9] projects data into subspaces to avoid performance degradation due to high-dimensionality and deploys half-space chains to obtain a density-based outlieriness score.

In this paper, we extend the Sparse Data Observers (SDO) outlier detection algorithm [10]. SDO builds a low-density model for a dataset by randomly sampling $k \in \mathbb{N}$ points from it, which are termed *observers*. To avoid using outliers as observers, the set of observers is cleaned based on a quality metric. An outlier score is obtained as the median distance to the $x \in \mathbb{N}$ closest observers.

Compared to the above methods for streaming data, SDOstream benefits from a model as AnyOut, but avoids the drawbacks of using clustering. It is not limited by a SW, so it avoids suffering from *short-sightedness* (Fig. 1). It operates in linear time like Random Forests, but, by keeping space locations in the model, its definition of outlieriness preserves data shapes instead of solely the trained model.

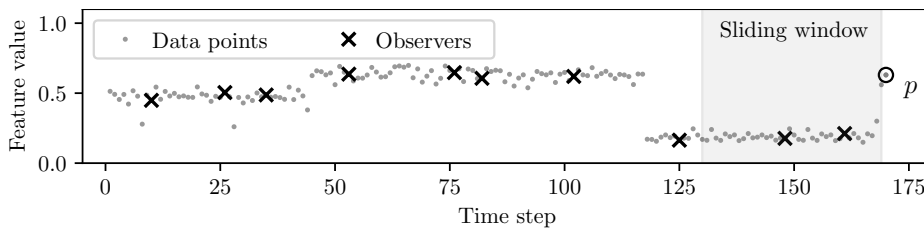


Fig. 1: p is an outlier for the SW, but an inlier for SDOstream.

3 SDOstream

The SDO algorithm can be adapted for streaming operation in a very natural way. We can adopt the usage of a fixed-size set of data points as observers, representing a model, consider a fixed-size fraction of observers as idle and compute an outlieriness score as median of distances to the x closest active observers. Streaming operation requires permanent adjustment of the model to new data. Hence, we introduce additional techniques to update observer quality metrics and observers themselves as new data arrives.

Table 1: Notation.

$f \in (0,1)$	Fading parameter.
$k \in \mathbb{N}$	Number of observers.
$x \in \mathbb{N}$	Number of neighbors to consider.
$q_{\text{id}} \in (0,1)$	Fraction of observers to consider idle.
Ω	Set of observers.
$\mathcal{N} \subset \Omega$	Set of x closest observers.
$\mathcal{N}_a \subset \Omega$	Set of x closest active observers.
$P_\omega \in \mathbb{R}^+$	Observations by ω .
$H_\omega \in \mathbb{R}^+$	Observer ω 's age.
$\bar{P}_\omega = P_\omega / (1 - f^{H_\omega})$	Average observations by ω .
$i_{\text{LAO}} \in \mathbb{Z}$	Index number of observer added most recently.

Parameters. The most important parameter of SDOstream is $f \in (0,1)$, the fading parameter. f controls how quickly the model is able to adjust to newly observed data clusters and, on the other hand, how stable the model is with respect to noise in the processed data. It is beneficial to write f as $f = \exp(-T^{-1})$ with a time parameter $T \in \mathbb{R}^+$, which in its function can be best understood as similar to the window size of SW approaches. Furthermore, $k \in \mathbb{N}$ denotes the model size, i.e. the number of observers to use in stationary operation. Intuitively, k should be increased for highly diverse data. Finally, $q_{\text{id}} \in (0,1)$ denotes the fraction of observers to ignore for outlier scoring (called *idle* observers) and $x \in \mathbb{N}$ is the number of nearest observers to consider. The reader is referred to [10] for a discussion of how to choose these parameters.

Observer idle-active split. In SDO, observations $P_\omega \in \mathbb{N}$ for an observer ω , i.e. the number of data points for which ω is contained in the nearest-observers set, constitutes a quality metric, which is used as basis for determining the set of active observers. To adjust P_ω to newly seen data, while being able to scale to arbitrary time scales, we deploy an exponential moving averaging approach. In particular, we keep track of an exponential moving average $P_\omega \in \mathbb{R}^+$ of observations by ω . Hence, for each processed data point, we set $P_\omega \leftarrow fP_\omega + 1$ if it belongs to the nearest-observers set of the current data point, and $P_\omega \leftarrow fP_\omega$ if it does not. We use P_ω to distinguish between active and idle observers, i.e., the q_{id} -fraction of observers with the lowest P_ω values are declared idle.

Replacing observers. To retain an up-to-date model, it is necessary to regularly replace outdated observers by new data points. It appears reasonable to select observers for removal based on the minimization of a quality metric like P_ω . However, new observers have a lower P_ω than established observers. While this is intended for the idle-active split, selecting observers for removal based on P_ω is ill-fated, as it would lead to new observers being replaced constantly. Instead, an average observer quality $\bar{P}_\omega \in \mathbb{R}^+$ is required, which sets P_ω in relation to the maximum value it can assume. Assume that data points arrive with a constant inter-arrival time (IAT) of $\delta \in \mathbb{R}^+$. Then, after a time $H \in \mathbb{R}^+$, i.e. after $\lfloor \frac{H}{\delta} \rfloor$ data points, the maximum P_ω can have assumed is $1 + f^\delta + \dots + f^{\lfloor \frac{H}{\delta} \rfloor \delta} = \frac{1 - f^{(\lfloor \frac{H}{\delta} \rfloor + 1)\delta}}{1 - f^\delta} \approx \frac{1 - f^H}{1 - f^\delta}$. Hence, $P_\omega / \frac{1 - f^H}{1 - f^\delta}$

can be used as a metric for selecting the observer to replace. Since we are only interested in comparing observers, the constant factor $1-\delta$ can be omitted, obtaining $\bar{P}_\omega = P_\omega/(1-f^H)$. Clearly, real-world IATs are not constant. However, it can be shown that the provided \bar{P}_ω remains valid also for non-constant IATs, given that the average IAT does not change at a time scale of model adjustment time.

Putting the blocks together. Algorithm 1 depicts SDOstream. Inputs consist of a data point $v_i \in \mathbb{R}^m$ and the corresponding time stamp $t_i \in \mathbb{R}$. In **line 1**, the x closest observers are determined from both sets of active and idle observers. Hence, the set of observers Ω is queried together with P_ω to determine the current idle-active split and find the sets \mathcal{N} and \mathcal{N}_a , containing the x closest observers and x closest active observers, respectively.

In **lines 2-3** observer quality P_ω and age H_ω of all observers are adjusted and, subsequently, the model itself. For this purpose, in **line 4**, the probability of sampling a data point to be used as observer, is proportional to $\frac{\sum_{\omega \in \mathcal{N}} P_\omega}{\sum_{\omega \in \Omega} P_\omega}$ to support an even P_ω and, hence, a representative model. Furthermore, note that during initial start-up using all available data points as observers might lead to a very poor model due to temporal dependence in the data stream, e.g. if the k first seen data points all belong to the same cluster. Hence, also during start-up, observers are added with a finite rate, which is gradually decreased to quickly approach the desired model size on the one hand, and sample over a large enough time period on the other hand. In stationary operation data points are then sampled with a rate of $k/T = -k \ln f$, i.e. within a time $T = -1/\ln f$ all observers are replaced one time on average.

Finally, in **line 11**, outlier scoring takes place and, similar to SDO, is performed by computing the median distance to the x closest active observers.

3.1 Time and Space Complexity

If n denotes the number of processed samples, we show that SDOstream has space and time complexity $O(k)$ and $O(n \log k)$, respectively. Indeed, considering solely dependence on n , space and time complexity are $O(1)$ and $O(n)$, as SDOstream operates with a fixed-size model. This benefit allows scaling to data streams of arbitrary size.

Algorithm 1 SDOstream: Processing a data point (v_i, t_i) .

- 1: Find x closest observer sets \mathcal{N} and \mathcal{N}_a
 - 2: Set $H_\omega \leftarrow H_\omega + (t_i - t_{i-1})$ and $P_\omega \leftarrow f^{t_i - t_{i-1}} P_\omega \quad \forall \omega \in \Omega$
 - 3: Set $P_\omega \leftarrow P_\omega + 1 \quad \forall \omega \in \mathcal{N}$
 - 4: **if** Ω empty **or** $r \leq -\ln(f) \frac{k^2}{x} \frac{\sum_{\omega \in \mathcal{N}} P_\omega}{\sum_{\omega \in \Omega} P_\omega} \frac{t_i - t_{i_{\text{LAO}}}}{i - i_{\text{LAO}}}$ with $r \in_R [0,1]$ **then**
 - 5: **if** $|\Omega| = k$ **then**
 - 6: Remove $\arg \min_{\omega \in \Omega} \bar{P}_\omega$ from Ω
 - 7: **end if**
 - 8: Add v_i to Ω
 - 9: Set $i_{\text{LAO}} \leftarrow i, P_{v_i} \leftarrow 1$ and $H_{v_i} \leftarrow 0$
 - 10: **end if**
 - 11: **return** $\text{median}_{\omega \in \mathcal{N}_a} (\|\omega - v_i\|)$ as outlier score
-

Furthermore, Algorithm 1 suggests $O(nk)$ run time when considering also model size k . However, Algorithm 1 mainly aims at simplicity of presentation and for highly dynamic, diverse streams, sub-linear run time with respect to k might be desired.

To obtain run time in $O(n \log k)$, observers can be stored in a data structure for efficient nearest neighbor search like an M-Tree [11], for which the authors observed logarithmic dependence of tree size. Tree-based structures can be used for keeping observers ordered with respect to P_ω . Hence, both nearest observer search and idle-active split are possible in $O(\log k)$. Discovery of the observer with lowest \bar{P}_ω takes $O(k)$ and is required every $\frac{T}{k\delta}$ th sample on average. Hence, SDOstream takes $O(n \log k + \frac{k\delta}{T}nk)$ amortized time. In the primary case of high-rate settings, i.e., $\delta/T \rightarrow 0$, total space and time complexity thus are $O(k)$ and $O(n \log k)$, respectively.

4 Performance Evaluation

We used real and artificial data to study SDOstream, using L2 distance and deploying genetic algorithms for tuning hyperparameters for all studied algorithms. To investigate the temporal behavior for transient events like suddenly emerging clusters, we crafted a synthetic dataset with 100,000 data points, 6 clusters and 2% outliers using MDCGen [12]. We removed data points for two clusters, so that the clusters appear after the 50,000th data point. On the right side, Fig. 2 shows the ensemble ROC-AUC and the average sampling rate computed over 10,000 runs of SDOstream with $T=10,000$ and appropriately randomly permuted datasets.

As shown, the clusters are first classified as outliers, but after 1,000 ($\frac{T}{10}$) to 2,000 ($\frac{2T}{10}$) data points, performance steeply increases, returning to stationary performance after about 4,000 ($\frac{4T}{10}$) data points. As a side effect of the factor $\frac{\sum_{\omega \in \mathcal{N}} P_\omega}{\sum_{\omega \in \Omega} P_\omega}$, the sampling rate is affected by the event to a very small extent.

On the left side of Fig. 2 we provide a comparison of performance results for the KDDCup'99 dataset, which contains 4.9 million data samples of streaming data. The window size is used as time parameter for xStream and SW-KNN, and T is used for SDOstream. To evade influence of the convergence phase, we only used the second half of returned outlierness scores for computing the ROC-AUC for both hyperparameter search and performance evaluation. As depicted, SDOstream exhibits very good performance in streaming scenarios, retains low processing time and shows to be stable over a wide range of different T .

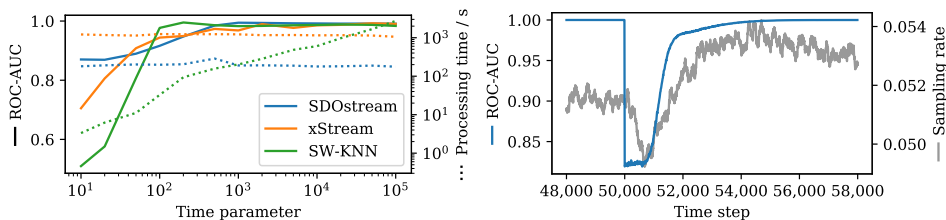


Fig. 2: Performances for the KDDCup'99 dataset (left) and behavior for emerging new clusters at $t=50,000$ with synthetic data (right).

Table 2: ROC-AUC scores for static datasets.

	SDOstream	SW-KNN	xStream	SDO	KNN	LOF	iForest
Annthroid	0.627	0.594	0.522	0.610	0.644	0.674	0.807
Cardiotocogr.	0.815	0.808	0.810	0.818	0.784	0.810	0.804
PageBlocks	0.904	0.903	0.923	0.910	0.904	0.943	0.921

Since only few datasets are established containing streaming data, and to provide a better comparison with popular algorithms, Table 2 compares SDOstream also against non-streaming algorithms using static datasets from [2], again evaluating stationary performance by considering the second half of returned scores. Here, SDOstream shows performances comparable to established outlier detection methods.

Concluding, SDOstream proved to be a fast, versatile outlier detection algorithm, demonstrating its strengths particularly when facing substantial data volumes by creating a representative model of the data. While providing an intuitive outlieriness definition, it is flexible and allows memory length to be set solely based on the application’s requirements.

References

- [1] M. M. Breunig, H.-P. Kriegel, et al. LOF: identifying density-based local outliers. In *Proc. of the 2000 ACM SIGMOD Int. Conf. on Management of data*, pages 93–104, 2000.
- [2] G. O. Campos, A. Zimek, et al. On the evaluation of unsupervised outlier detection: measures, datasets, and an empirical study. *Data Mining and Knowledge Discovery*, 30(4):891–927, 2016.
- [3] F. Angiulli and F. Fassetti. Detecting distance-based outliers in streams of data. In *Proc. of the 16th ACM Conf. on Information and Knowledge Management, CIKM’07*, pages 811–820, New York, NY, USA, 2007. ACM.
- [4] D. Yang, E. Rundensteiner, et al. Neighbor-based pattern detection for windows over streaming data. In *Proc. of the 12th Int. Conf. on Extending Database Tech.: Advances in Database Tech., EDBT’09*, pages 529–540, New York, NY, USA, 2009. ACM.
- [5] M. Kontaki, A. Gounaris, et al. Continuous monitoring of distance-based outliers over data streams. In *IEEE 27th Int. Conf. on Data Engineering*, pages 135–146, 2011.
- [6] I. Assent, P. Kranen, et al. Anyout: Anytime outlier detection on streaming data. In *Database Systems for Advanced Applications*, pages 228–242. Springer, 2012.
- [7] K. Wu, K. Zhang, et al. RS-Forest: A Rapid Density Estimator for Streaming Anomaly Detection. In *IEEE Int. Conf. on Data Mining*, pages 600–609, 2014.
- [8] S. Guha, N. Mishra, et al. Robust random cut forest based anomaly detection on streams. In *Proc. of The 33rd Int. Conf. on Machine Learning*, volume 48 of *Proc. of Machine Learning Research*, pages 2712–2721, New York, USA, 2016. PMLR.
- [9] E. Manzoor, H. Lamba, et al. xStream: Outlier detection in feature-evolving data streams. In *24th ACM SIGKDD Int. Conf. on Know. Discovery and Data Mining*, 2018.
- [10] F. Iglesias, T. Zseby, et al. Outlier detection based on low density models. In *2018 IEEE Int. Conf. on Data Mining Workshops (ICDMW)*, pages 970–979, 2018.
- [11] P. Ciaccia, M. Patella, et al. M-tree: An efficient access method for similarity search in metric spaces. In *Proc. of the 23rd VLDB conference*, pages 426–435, 1997.
- [12] F. Iglesias, T. Zseby, et al. MDCGen: Multidimensional Dataset Generator for Clustering. *Journal of Classification*, 36:599–618, 2019.