

# A comprehensive analysis of concept drift locality in data streams

Gabriel J. Aguiar<sup>a</sup>, Alberto Cano<sup>a</sup>

<sup>a</sup>*Department of Computer Science, Virginia Commonwealth University,  
Richmond, Virginia, USA*

---

## Abstract

Adapting to drifting data streams is a significant challenge in online learning. Concept drift must be detected for effective model adaptation to evolving data properties. Concept drift can impact the data distribution entirely or partially, which makes it difficult for drift detectors to accurately identify the concept drift. Despite the numerous concept drift detectors in the literature, standardized procedures and benchmarks for comprehensive evaluation considering the locality of the drift are lacking. We present a novel categorization of concept drift based on its locality and scale. A systematic approach leads to a set of 2,760 benchmark problems, reflecting various difficulty levels following our proposed categorization. We conduct a comparative assessment of 9 state-of-the-art drift detectors across diverse difficulties, highlighting their strengths and weaknesses for future research. We examine how drift locality influences the classifier performance and propose strategies for different drift categories to minimize the recovery time. Lastly, we provide lessons learned and recommendations for future concept drift research. Our benchmark data streams and experiments are publicly available at <https://github.com/gabrieljaguiar/locality-concept-drift>.

*Keywords:* Data streams, Multi-Class, Concept Drift, Machine Learning

---

## 1. Introduction

Modern data sources continuously generate information characterized by both volume and velocity, flooding learning systems with a constant flow of data. This scenario is commonly referred to as data streams [1, 2]. Traditional classification methods, designed for static data, struggle to keep

up with the ever-changing characteristics of these incoming instances [1, 3]. Given the dynamic nature of data streams, it becomes essential for learning methods to adapt and acquire knowledge about emerging concepts over time. This phenomenon is known as concept drift [4], and it can manifest in various ways, including shifts in class distribution and decision boundaries [5], and the emergence of new features or classes [6]. If not detected and addressed effectively, concept drift can significantly degrade predictive performance, as knowledge learned from older concepts may not be useful anymore to classify recent instances [7].

In recent years, the issue of concept drift has garnered significant attention within the research community across various domains, including sensors, robotics, system monitoring, and anomaly detection [8]. Current research in this field is tackling increasingly complex challenges. These challenges include accurately detecting concept drift within unstructured and noisy datasets [9], providing understandable explanations for concept drift [10], and effectively responding to drift by adapting relevant knowledge [11]. When we extend these concerns to scenarios involving multiple classes, we encounter a complex and perplexing scenario that actually occurs in many real-life applications. Detecting concept drift in such contexts becomes exceptionally demanding, as we must account for the evolving nature of multiple classes [6, 12]. In addition to the challenges previously mentioned, it is important to note that the location of concept drift within the feature space significantly influences both the performance of classifiers and the effectiveness of drift detection methods [6, 13]. However, there is a lack of studies that evaluate drift detectors under varying degrees of drift locality or provide benchmark datasets to support research in this crucial area.

**Motivation.** While the literature offers numerous concept drift detectors, there remains a notable absence of standardized procedures and benchmarks for a comprehensive assessment of these methods when considering the locality of concept drifts. Specifically, there is a lack of dedicated benchmarks suitable for evaluating drift detectors across a diverse spectrum of challenges, particularly those tied to the locality of concept drift. An in-depth experimental comparison of state-of-the-art drift detectors, applied to a diverse set of challenges, would provide valuable insights into the performance of these detectors under various conditions. Moreover, existing studies often focus on specific subsets of detectors and data challenges, typically limited to binary class data. These studies often fail to provide insights into which aspects of concept drift should be considered. Therefore, we propose a comprehen-

sive study to understand and assess the performance of drift detectors across a wide array of difficulties and analyze how the locality of a concept drift impacts its detection.

**Overview and main contributions.** This paper presents a comprehensive study for benchmarking and evaluating the impact of the locality and magnitudes of a concept drift on the classifier or drift detector. We systematically identify critical challenges within this domain and leverage them to create a set of benchmark problems that encompass various difficulties, guided by a novel concept drift categorization. Furthermore, we conduct a comparative evaluation of nine state-of-the-art drift detectors across this wide range of difficulty. This analysis not only identifies the top-performing detectors but also sheds light on their specific strengths and weaknesses, providing valuable insights for future research in drift detection. The main contributions of this paper can be summarized as follows:

- **Concept drift locality categorization.** We present a novel categorization that organizes concept drift based on the number of affected classes and the magnitude of the change. This categorization empowers us and future researchers to evaluate the effect of various concept drift levels when proposing new drift detectors or classifiers. This categorization guides the creation of benchmark problems, ranging from scenarios affecting only one class to those that transform the entire data stream. The comprehensive set of benchmark problems consists of 2,760 data stream benchmarks.
- **Drift locality impact evaluation.** We present a comprehensive study designed to assess the influence of concept drift locality on its detection, encompassing both binary and multi-class data streams.
- **Comparison between state-of-the-art drift detectors.** We conduct an extensive, comprehensive, and reproducible comparative study among 9 state-of-the-art drift detectors based on the proposed framework and 2,760 benchmark problems.
- **Recommendations and open challenges.** Based on the results from the experimental study, we derive recommendations seeking insights into the strengths and weaknesses of the top-performing drift detectors. These recommendations aim to provide a comprehensive understanding of the detectors' capabilities. Additionally, we identify

open challenges within the domain of learning from data streams impacted by concept drift, outlining directions for future research.

This paper is structured as follows: Section 2 provides the theoretical foundation and discusses related work. Section 3 introduces our proposed concept drift locality categorization and presents the benchmark problems. Section 4 outlines the experimental setup, while Section 5 presents the results. Section 6 delves into the lessons learned, and Section 7 offers the conclusion and outlines directions for future work.

## 2. Background and related work

This section reviews the background and related work. We provide an overview of the literature on data streams and concept drift detection.

### 2.1. Data streams

A data stream refers to a potentially unbounded sequence of ordered instances that arrive over time within a system. Learning from data streams imposes specific limitations on classifiers [3]. We can define a stream, denoted as  $S$ , as a sequence  $\langle s_1, s_2, s_3, \dots, s_\infty \rangle$ , where  $s_i = (X, y)$ . This stream can be handled either one instance at a time (online scenario) or in batches (block scenario). Data streams exhibits four primary characteristics [2, 5]: (i) Volume, (ii) Velocity, (iii) Veracity, and (iv) Non-stationarity, which present challenges to classifiers that must adapt accordingly.

Data streams are susceptible to a phenomenon known as concept drift [14, 15]. Each instance arrives at a specific time, denoted as  $t$ , and is generated based on a probabilistic distribution denoted as  $\phi^t(X, y)$ , where  $X$  represents the feature vector and  $y$  denotes the class label. If all instances in the stream are generated by the same probability distribution, the data is considered stationary, indicating a single and stable underlying concept. However, in real-world applications, data rarely adheres to stationary assumptions [16]. On the other hand, when two distinct instances arriving at times  $t$  and  $t+C$  are generated by  $\phi^t(X, y)$  and  $\phi^{t+C}(X, y)$  respectively, and if  $\phi^t \neq \phi^{t+C}$ , indicates the occurrence of a concept drift. This phenomenon impacts various aspects of a data stream and, as such can be examined from multiple viewpoints. When analyzing and comprehending concept drift, the following factors come into consideration [6, 5]:

- **Influence of the decision boundaries.** Firstly, it is necessary to consider how concept drift affects the learned decision boundaries, distinguishing between real and virtual concept drifts. Virtual drift produces a change in the unconditional probability distribution  $P(x)$ , without affecting the learned decision boundaries. Although virtual drift does not impair learning models, its detection is necessary to avoid false alarms and prevent unnecessary, costly adaptations. In contrast, real concept drift modifies the decision boundaries, making them worthless to the current concept. Detecting and adapting to real concept drift is crucial for preserving predictive performance.
- **Speed of changes.** There are three types of concept drift [17]: (i) incremental; (ii) gradual; and (iii) sudden concept drifts as illustrated in Fig. 1. Incremental drift generates a sequence of intermediate states between the old and new concept, while gradual drift oscillates between instances coming from both old and new concepts, with the new concept becoming more and more frequent over time. Finally, sudden drift instantaneously switches between old and new concepts, leading to an instant degradation of the underlying learning algorithm.
- **Recurrence.** Changes in the stream can be either unique or recurring. In the latter case, the previously seen concept may reemerge over time, allowing us to recycle previously learned knowledge. The past knowledge can be used as an initialization point for the drift recovery.
- **Presence of noise.** Noise can take the form of sporadic, insignificant variations within a stream that can be disregarded, or substantial corruption within the features or class labels that need to be dealt with in order to prevent the input of misleading or adversarial data into the classifier [18]. Drift detectors must ignore noise.
- **Locality.** The literature distinguishes between global and local concept drifts [13]. The former impacts the entire data stream, while the latter pertains to specific parts of it, such as individual clusters of instances or subsets of classes. However, this classification is too shallow and a more in-depth and detailed discussion regarding the concept drift’s locality will be introduced in this manuscript.

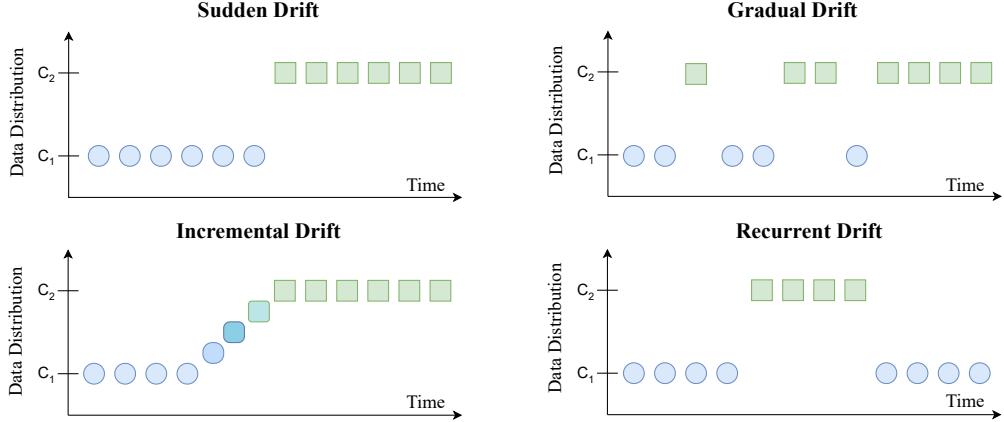


Figure 1: Different types of drifts depending on their speed.

To address the challenges posed by concept drift, two approaches are commonly employed: (i) implicit and (ii) explicit. Implicit methods manage drift adaptation through intrinsic mechanisms integrated within the classifier, assuming its capability to self-adjust to new instances reflecting the most recent concept while gradually discarding outdated information [19, 8]. These approaches involve establishing appropriate learning and forgetting rates, utilizing adaptive sliding windows, or continually tuning hyperparameters. Conversely, explicit approaches assign drift adaptation to an external tool known as a drift detector [8, 15]. Drift detectors continuously monitor stream properties (e.g., statistics) or classifier performance (e.g., error rates). They raise a warning signal when there are indications of impending drift and trigger an alarm signal when concept drift has occurred.

## 2.2. Drift detectors

In recent years, a plethora of explicit drift detectors have been introduced [20]. As defined by Lu et al. [15], drift detectors function by extracting critical features from both historical and newly arrived data and subsequently subjecting them to dissimilarity tests. Like classifiers, drift detection methods fall into three categories: supervised (or error rate-based), semi-supervised, and unsupervised (or data distribution-based)[15]. The primary distinction among these categories is the point at which drift is identified. While supervised drift detectors pinpoint alterations in class boundaries, unsupervised detectors focus on tracking shifts in data distribution [4].

Among supervised drift detectors, the most widely adopted group relies on evaluating classifier error or accuracy using labeled instances. This category of drift detectors can be further divided into three distinct types: (i) change detection-based detectors, (ii) statistical-based detectors, and (iii) window-based detectors.

- **Change detection-based detectors.** In the first category, we find detectors like Page Hinkley [21], CUSUM [21], and Geometric Moving Average [22]. These detectors utilize cumulative sums to trigger an alarm when a significant change in input data occurs. While they are computationally efficient, their performance heavily depends on the selection of hyperparameters.
- **Statistical-based detectors.** The second group comprises detectors such as the Drift Detector Method (DDM) [1] and its variants like Early DDM (EDDM) [23], Reactive DDM (RDDM) [11], DDM based on Hoeffding’s bound (HDDM) [24], Wilcoxon Rank Sum Test Drift Detector (WSTD) [25], and EWMA for Drift Detection (ECDD) [26]. These methods compute statistical features over time based on error rates. They operate under the assumption that as long as the data distribution remains stationary, the learner’s error rate will decrease with an increasing number of analyzed samples. When the error rate exceeds a predefined threshold, an alert is triggered.
- **Window-based detectors.** Finally, a popular category of supervised detectors relies on metrics calculated within subwindows of a data stream. An example is the ADaptive WINdow (ADWIN) [27], which employs an adaptive sliding window based on Hoeffding’s inequality. ADWIN manages a sliding window divided into two sub-windows representing old and new data and dynamically adjusts its size, expanding during periods of stability and contracting in the presence of drift. ADWIN signals a drift when the average between the two windows surpasses a predefined threshold. This approach has inspired the development of several other detectors such as Kolmogorov-Smirnov Windowing (KSWIN) [28], Statistical Test of Equal Proportions (STEPD) [29], and others [30, 31].

In addition to detectors based on classifier accuracy, there are also supervised trainable detectors that employ traditional machine learning classifiers to detect concept drift. Examples of such detectors include the Restricted Boltzmann Machine (RBM-IM) [6], the Complexity Drift Detector (C2D) [32] and QuadCDD [33]. It is important to note that there is a notable absence of drift detectors specifically designed for multi-class scenarios, which come with their own unique characteristics and challenges.

Unsupervised drift detectors are primarily focused on identifying disparities in unlabeled data without the need for added supervision. These detectors typically employ a distance metric to quantify the dissimilarity between the distribution of historical data and newly arrived data [34, 35]. When this dissimilarity is statistically significant, it triggers a process to update the learning model. In this category, we find detectors such as Statistical Change Detection for multi-dimensional data (SCD)[36], the PCA-based change detection framework (PCA-CD)[37], Equal Density Estimation (EDE)[38], Least Squares Density Difference-based Change Detection Test (LSDD-CDT)[39], among others [15]. More advanced unsupervised methods aim to precisely locate where the drift occurred in the feature space. These detectors focus on spatial searches employing various dissimilarity measures [10, 40]. They tackle concept drift at its root, addressing distribution drift. Typically, these algorithms require users to define both a historical time window and a new data window. A common approach is to use two sliding windows, with the historical time window fixed while sliding the new data window [41].

As aforementioned, numerous drift detectors have been proposed in the literature, highlighting the need for comparisons to understand how each one behaves in different complex scenarios. While Lu et al. [15] and Suárez-Cetrulo et al. [8] reviewed the state-of-the-art in data stream learning with a focus on drift detectors, however they lacked an empirical comparison between these detectors. Additionally, the locality of the concept drift is often neglected [6, 13]. In a different approach, Barros and Santos [20] conducted an extensive empirical comparison and evaluation of concept drift detection methods across various data stream configurations. Moreover, Table 1 presents a summary of works that empirically compare drift detectors and the type of concept drift evaluated. However, a notable gap exists in the literature concerning concept drift detection in data streams with multiple classes and how the locality of the drift influences its detection.

Table 1: Comparison of most commonly applied, drift detectors and types of concept drift by contributions in related works

Reference	Drift Detectors	Drift Type
Gama et al. [1]	DDM	Sudden; Gradual
Baena-Garcia et al. [23]	DDM, EDDM	Sudden; Gradual
Frias-Blanco et al. [24]	DDM, ADWIN, HDDM <sub>W</sub> , ECDD	Sudden; Gradual
Gonçalves Jr et al. [42]	DDM, EDDM, ADWIN, ECDD, STEPD, PH, NB, PL, DOF	Sudden; Gradual
Barros and Santos [20]	DDM, EDDM, ADWIN, HDDM <sub>W</sub> , HDDM <sub>A</sub> , ECDD, SEQDRIFT, SEED, STEPD, FHDDM, FTDD, RDDM, WSTD	Sudden; Gradual
Santos et al. [43]	DDM, EDDM, ADWIN, HDDM <sub>W</sub> , HDDM <sub>A</sub> , ECDD, SEQDRIFT, STEPD, FHDDM, RDDM, and WSTD	Sudden; Gradual
Babüroğlu et al. [44]	DDM, EDDM, ADWIN, HDMM <sub>W</sub> , HDDM <sub>A</sub> , ECDD, SEQDRIFT, SEED, STEPD, FHDDM, FTDD, RDDM, WSTD, GMA, PH	Sudden
Korycki and Krawczyk [6]	FHDDM, RDDM, WSTD, PerfSim, DDM-OCI, RBM-IM	Sudden; Gradual; Incremental
Poenaru-Olaru et al. [45]	DDM, EDDM, ADWIN, HDDM <sub>W</sub> , HDDM <sub>A</sub>	Sudden; Gradual
Sakurai et al. [46]	DDM, EDDM, HDDM <sub>W</sub> , HDDM <sub>A</sub>	Sudden; Gradual; Incremental

### 3. Proposed categorization of concept drift locality

As discussed in various studies [6, 13, 47, 48], local data difficulty factors and the number of classes exert a substantial influence on the capabilities of classifiers when handling data streams. Furthermore, these factors also play a major role in concept drift detection. Notably, the performance of standard drift detectors can vary significantly based on the underlying distribution of the stream, especially when they lack strategies to address these factors. This can result in either excellent performance or suboptimal outcomes.

However, the current body of literature addressing the categorization and characterization of concept drifts in streaming classification problems often overlooks two crucial factors: the number of classes affected and local data difficulties. It is important to note that a significant portion of concept drift research primarily focuses on binary classification and global shifts in data distribution, or perturbations that primarily impact the minority class, as observed in the imbalanced scenario [47]. Given the substantial impact that local data difficulty factors and the number of classes affected can have on non-stationary data streams, the existing categorizations of concept drift may fall short in evaluating concept drift detection in dynamic environments. They may fail to encompass all the nuances that make a concept drift particularly challenging or solvable. Therefore, we propose an extended concept drift categorization that explicitly incorporates considerations of locality and the number of classes affected when assessing concept drift. This expanded framework should pave the way for further research in the development of classifiers and drift detectors tailored for multi-class data streams in non-stationary environments, enabling systematic studies under different relevant drifting conditions.

#### 3.1. Categorization of concept drift locality

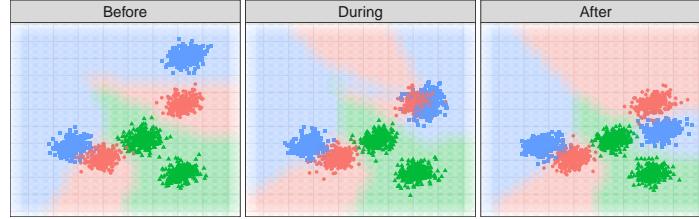
First, we formally define the proposed categorization. Consider a bounded  $d$ -dimensional attribute space  $X$  and output space  $y$  and a given posterior probability distribution  $\phi^t(X, y)$  at time  $t$ , which can be expanded as  $\phi_{c_1}^t(X, y) \cup \phi_{c_2}^t(X, y) \dots \cup \phi_{c_n}^t(X, y)$  where  $c_i$  represents the distribution for a given class and  $n$  the number of classes. As previously defined if  $\phi^t(X, y) \neq \phi^{t+C}(X, y)$  a concept drift has occurred and falls into one of the four categories:

- **Single-Class Local Concept Drift:** Given the drifted distribution  $\phi^{t+C}(X, y)$  exists **only one**  $i \in [0, n]$  that  $\phi_{c_i}^{t+C}(X, y) = \phi_{c_i}^t(X, y) \cup \hat{\phi}_{c_i}(X, y)$ .
- **Single-Class Global Concept Drift:** Given the drifted distribution  $\phi^{t+C}(X, y)$  exists **only one**  $i \in [0, n]$  that  $\phi_{c_i}^{t+C}(X, y) = \hat{\phi}_{c_i}(X, y)$ .
- **Multi-Class Local Concept Drift:** Given the drifted distribution  $\phi^{t+C}(X, y)$  exists **more than one**  $i \in [0, n]$  that  $\phi_{c_i}^{t+C}(X, y) = \phi_{c_i}^t(X, y) \cup \hat{\phi}_{c_i}(X, y)$
- **Multi-Class Global Concept Drift:** Given the drifted distribution  $\phi^{t+C}(X, y)$  exists **more than one**  $i \in [0, n]$  that  $\phi_{c_i}^{t+C}(X, y) = \hat{\phi}_{c_i}(X, y)$ .

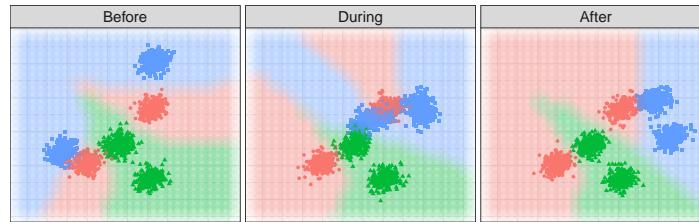
Figs. 2 and 3 present a visual representation of the four defined categories. Local drifts occur when the original concept and a new distribution emerge without overlapping. This type of concept drift can be exceedingly subtle yet has the potential to compromise predictive accuracy in specific regions of the feature space. Consequently, it poses a challenge for detection and resolution. On the other hand, global concept drift involves a complete transformation of the data distribution, often resulting in less subtle changes. However, the ease of drift detection depends on the behavior of the new distribution. It is worth noting that learning mechanisms and concept drift detection strategies can differ significantly in scenarios with multiple classes [6, 48]. Thus, when assessing classifier performance or concept drift detection, we must also account for the number of classes affected by the concept drift. This proposed categorization enables us to analyze concept drift detection while considering intricacies that have not received extensive attention in the existing literature. Additionally, it empowers the development of novel classifiers and drift detectors capable of handling a spectrum of concept drifts, ranging from subtle to noticeable.

### 3.2. Benchmarks

To explicitly assess the performance of classifiers and drift detectors in data streams featuring the concept drift categories outlined earlier, we introduce a set of drift difficulties corresponding to each category within our proposed framework. These difficulties were implemented using two widely recognized data stream generators, namely Random RBF and Random Tree,

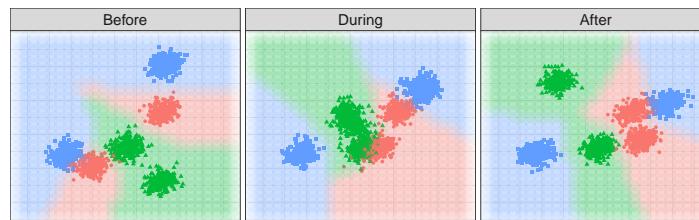


(a) Single-Class Local

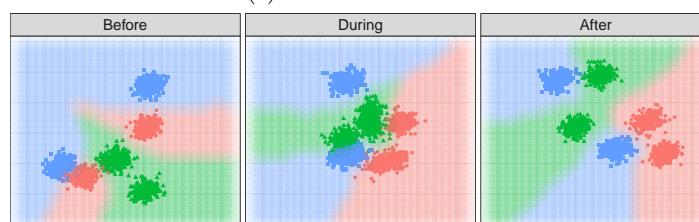


(b) Single-Class Global

Figure 2: Illustrative data distribution before, during and after a Single-Class Local (a) and Global (b) concept drift. Each color represents one class and background color represents decision boundaries.



(a) Multi-Class Local



(b) Multi-Class Global

Figure 3: Illustrative data distribution before, during and after a Multi-Class Local (a) and Global (b) concept drift. Each color represents one class and background color represents decision boundaries.

as established in the literature [5]. Table 2 presents a comprehensive overview of the proposed data difficulties along with their respective categorizations. For each difficulty setting, we generated data streams encompassing a range of class counts, including  $\{2, 3, 5, 10\}$  classes, and varying feature dimensions of  $\{2, 5, 10\}$ . Additionally, we introduced three distinct types of drifts: Sudden, Gradual, and Incremental. In the context of Multi-Class drifts, we considered scenarios where drifts affected subsets of classes, including  $\{\{2\}, \{2, 3\}, \{2, 3, 5\}, \{2, 3, 5, 10\}\}$ , depending on the number of classes.

Finally, to ensure reproducibility and facilitate the utilization of our proposed categorization and benchmarks, all the data streams used in our experiments have been made publicly available for future research.<sup>1</sup>

#### 4. Experimental Setup

The experiments are designed to assess the performance of drift detectors across a diverse range of drift challenges. We also investigate how the locality and the number of affected classes influence drift detection. The primary

---

<sup>1</sup><https://github.com/gabrieljaguaiar/locality-concept-drift>

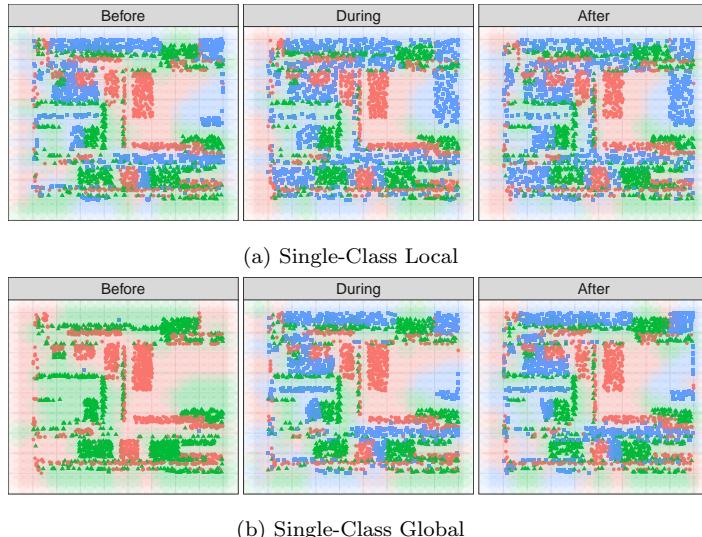


Figure 4: Data distribution before, during and after a Single-Class Local (a) and Global (b) `emerging_branch` concept drift. Each color represents one class and background color represents decision boundaries.

Table 2: Concept drift locality benchmark specifications. S: Sudden, G: Gradual, I: Incremental.

Generator	Impact on the class space	Impact on the feature space	Difficulty	Description	Drift Speed
RBF	Single-Class	emerging_cluster	A new subcluster emerge	S/G	
		reappearing_cluster	Part of the clusters of ONE class disappears for a period of time	S/G	
		splitting_cluster	Part of the clusters of ONE class splits into two clusters that move to another direction	S/G/I	
		merging_cluster	Part of the clusters of ONE class clusters are merged in the midpoint between them	S/G/I	
		moving_cluster	Part of the clusters of ONE class center moves to another position	S/G/I	
	Multi-Class	reappearing_cluster	All clusters of ONE class disappear for a period of time	S/G	
		splitting_subcluster	All subclusters of ONE class splits into two subclusters that move to another direction	S/G/I	
		merging_cluster	All subclusters of ONE class of one class merge	S/G/I	
		moving_cluster	All clusters of ONE class moves to another position	S/G/I	
		class_emergence	A new class appear	S/G	
RT	Single-Class	emerging_cluster	A new emerging subcluster for N classes appears	S/G	
		reappearing_cluster	Part of the clusters of N classes disappears for a period of time	S/G	
		splitting_cluster	Part of the clusters of N classes splits into two clusters that move to another direction	S/G/I	
		merging_clusters	Part of the clusters clusters are merged in the midpoint between them for N classes	S/G/I	
		moving_cluster	Part of the clusters of N classes center moves to another position	S/G/I	
	Multi-Class	swap_cluster	Part of the clusters of N classes swap position	S/G	
		reappearing_cluster	All clusters of N classes disappear for a period of time	S/G	
		splitting_cluster	All clusters of N classes splits into two clusters that move to another direction	S/G/I	
		merging_cluster	All clusters of N classes center merges	S/G/I	
		moving_cluster	All clusters of N classes moves to another position	S/G/I	
		swap_cluster	All clusters of N classes swaps position	S/G	
SVM	Local	emerging_branch	A new branch appears	S/G	
		prune_regrowth_branch	Part of the branches of ONE class are pruned and then regrowth this branches	S/G	
		prune_growth_new_branch	Part of the branches of ONE class are pruned and then other branches grow	S/G	
		prune_regrowth_branch	All of the branches of ONE class are pruned and then regrowth this branches	S/G	
		prune_growth_new_branch	All of the branches of ONE class are pruned and then other branches grow	S/G	
	Global	class_emerging	A new class appear	S/G	
		emerging_branch	A new branch appears for N classes	S/G	
		prune_regrowth_branch	Part of the branches of N classes are pruned and then regrowth this branches	S/G	
		prune_growth_new_branch	Part of the branches of N classes are pruned and then other branches grow	S/G	
		split_node	Part of the leaves of N classes are split to generate two leafs of two different classes	S/G	
NB	Global	swap_leaves	Part of the leaves of N classes swaps with another class	S/G	
		prune_regrowth_branch	All of the branches of N classes are pruned and then regrowth this branches	S/G	
		prune_growth_new_branch	All of the branches of N classes are pruned and then other branches grow	S/G	
		split_node	All of the leaves of N classes are split to generate two leafs of two different classes	S/G	
		swap_leaves	All of the leaves of N classes swaps with another class	S/G	

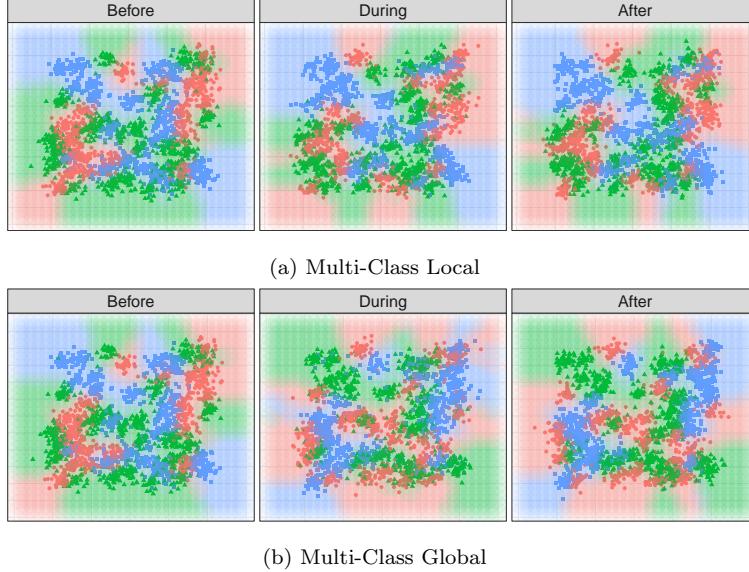


Figure 5: Data distribution before, during and after a Multi-Class Local (a) and Global (b) `swap_cluster` concept drift with all classes affected. Each color represents one class and background color represents decision boundaries.

goal is to gain insights into the performance of various drift detectors and how each type of drift challenge affects their effectiveness. Additionally, we examine the impact of each challenge on classifiers that rely on explicit drift detectors. To address these objectives, we formulate the following research questions (RQ):

- **RQ1:** Which concept drift detector demonstrates the most effective detection performance across all scenarios?
- **RQ2:** Which concept drift detector excels in detection performance for each distinct type of scenario?
- **RQ3:** How does each category of concept drift influence the performance of the drift detector?
- **RQ4:** Which specific difficulty presents the most challenging scenario?
- **RQ5:** How does the number of classes and features affect the concept drift detection?

- **RQ6:** How does each scenario impact the performance of classifiers utilizing the best-performing drift detector?

#### 4.1. Drift detectors

To assess the performance of drift detectors across various types of concept drift, we chose 9 state-of-the-art supervised drift detectors displayed in Table 3. We made these selections based on their widespread use in the literature and their demonstrated effectiveness in numerous scenarios. It is important to note that we included at least one representative from each group of supervised drift detectors mentioned in Section 2. This choice was deliberate, aiming to gain insights into the strengths and weaknesses of each group.

To monitor drift, we leverage the error distribution of the classifier to effectively monitor concept drifts, as suggested in the existing concept drift literature [27]. The error distribution is binary, and we determine the value to update the drift detector using the equation:

$$\phi = \begin{cases} 1, & \text{if } L(x) = y \\ 0, & \text{if } L(x) \neq y \end{cases} \quad (1)$$

where  $\phi$  represents the error distribution,  $x$  represents a new instance,  $y$  denotes the correct class, and  $L$  the classifier. By employing this approach, we can effectively detect concept drifts when there are changes in the classifier error distribution.

Table 3: Drift detectors used in the experiments.

Acronym	Name	Reference
ADWIN	ADaptive Windowing	Bifet and Gavalda [27]
KSWIN	Kolmogorov-Smirnov Windowing	Raab et al. [28]
STEPD	Statistical Test of Equal Proportions	Nishida and Yamauchi [29]
DDM	Drift Detection Method	Gama et al. [1]
EDDM	Early Drift Detection Method	Baena-Garcia et al. [23]
HDDM	Drift Detection Method based on Hoefing bound	Frias-Blanco et al. [24]
RDDM	Reactive Drift Detection Method	Barros et al. [11]
ECDD	EWMA for Drift Detection	Ross et al. [26]
PH	Page Hinkley	Page [21]

#### 4.2. Classifier

To evaluate the performance of the drift detectors, we opted to use Hoeffding Tree (HT) [49] as our classifier. This classifier is known for its exceptional predictive performance and lacks inherent drift detection or adaptation mechanisms, ensuring an unbiased assessment of concept drift detection.

#### 4.3. Performance evaluation

Given that each data stream in our experiments is characterized by a single known drift point, we employed an evaluation process similar to that used in classification tasks. Consequently, we extracted True Positive (TP), False Positive (FP), and False Negative (FN) values according to the following definitions:

- **TP:** When a concept drift detected at time  $t$  and  $t \in [i_d, i_d + R]$  where  $i_d$  is the instance where the concept drift happened, and  $R$  is the range of detection.
- **TN:** When no alert is raised and there is no known concept drift.
- **FP:** When a concept drift alert is raised at  $t$  and  $t \notin [i_d, i_d + R]$  where  $i_d$  is the instance where the concept drift happened, and  $R$  is the range of detection..
- **FN:** When no alert is raised within the range of detection  $R$ .

Subsequently, we calculated three performance metrics: Precision (Eq. 2), Recall (Eq. 3), and F1-Score (Eq. 4). Precision assesses the ratio between false alarms and accurate drift alerts, while recall measures the ratio between correct drift detection and undetected drifts. The F1-Score combines both of these metrics. These three metrics in combination provide insights into the behavior of each concept drift detector.

$$Precision = \frac{TP}{TP + FP} \quad (2) \qquad Recall = \frac{TP}{TP + FN} \quad (3)$$

$$F1 - Score = \frac{2 * (Precision * Recall)}{Precision + Recall} \quad (4)$$

Furthermore, for each accurate drift alarm, we calculate the average detection delay as described in Eq. 5, where  $D$  represents the instance position

of the drift, and  $D_d$  indicates the instance position when the drift alarm was triggered. This metric assesses the speed of drift detection.

$$Delay = \frac{\sum(D_d - D)}{TP} \quad (5)$$

#### 4.4. Reproducibility

The source code of the algorithms, the experiments and the data streams are publicly available on GitHub to facilitate the transparency and reproducibility of this research <sup>2</sup>. All results, interactive plots and tables are available on the website <sup>3</sup>. All the experiments, generators and drift detectors were implemented using Python 3.8 and the `river` [50] package. Algorithms were run on a GNU/Linux cluster with 192 Intel Xeon cores, 6 TB RAM, and Centos 7.

## 5. Results

This section presents and discusses the experimental results. Firstly, we conducted a comparative analysis of all drift detectors, irrespective of the specific scenarios, to determine which drift detector exhibits the best overall performance both with and without concept drift. Secondly, we investigate deeper into each scenario to gain insights into how the locality and the proposed data difficulties influence drift detection. Additionally, we evaluated the influence of the number of classes on drift detection in each scenario. Finally, we assessed how these proposed difficulties affect the classifier's performance and the impact of using a classifier in combination with an explicit drift detector.

### 5.1. Overall aggregated comparison for all scenarios

To address **RQ1**, we first evaluated the performance of all classifiers with a focus on data streams that either had or did not have concept drift, irrespective of how the concept drift affected the feature space or class space. Table 4 presents the results for data streams without concept drift. In these cases, as the stream remains stationary, we display only the average True

---

<sup>2</sup><https://github.com/gabrieljaguaiar/locality-concept-drift>

<sup>3</sup>Interactive plots and tables for all experiments are available at <https://gabrieljaguaiar.github.io/comprehensive-concept-drift/>

Negatives and False Positives values. Notably, only 3 drift detectors did not raise any alert in this scenario, highlighting the sensitivity of explicit drift detectors to changes in the error distribution, even when data distribution remains stable. Moreover, ECDD, EDDM, and STEPD exhibited a high level of sensitivity to any change, resulting in a high number of False Positives. In contrast, DDM and ADWIN displayed the lowest values of false alerts, demonstrating their stability over time.

Table 5 presents the results considering data streams with concept drift. It is interesting to note that EDDM exhibited the lowest delay among all evaluated drift detectors and the second-highest recall, although this came at the cost of raising numerous drift alerts, leading to a precision of 0%. A similar pattern is observed for ECDD and STEPD. When comparing DDM and its variations, their detection performance is quite similar, with HDDM displaying the best F1-Score and the shortest detection delay, highlighting the efficiency of Hoeffding’s inequality. In terms of F1-Score, ADWIN and PH displayed the best detection performance, with PH having a shorter detection delay.

In summary, ADWIN, DDM, and PH stood out as strong performers in both evaluated scenarios. They exhibited caution when dealing with stationary streams and demonstrated good detection performance when concept drift occurred. On the other hand, EDDM, ECDD, and STEPD could detect concept drifts quickly but raised a significant number of false alarms in the process.

Table 4: Comparison between drift detectors considering scenarios without concept drift.

Drift Detector	TN	FP
ADWIN	0.00	3.79
DDM	0.17	<b>2.71</b>
ECDD	<b>0.46</b>	904.58
EDDM	0.00	62914.04
HDDM	0.00	41.83
KSWIN	0.08	10.25
PH	0.00	12.46
RDDM	0.00	27.96
STEPD	0.00	4628.13

Table 5: Comparison between drift detectors considering all evaluated drift difficulties.

Drift Detector	Precision	Recall	F-1	Delay
ADWIN	<b>7.40%</b>	34.50%	<b>12.19%</b>	2192
DDM	1.82%	5.45%	2.73%	2168
ECDD	0.03%	36.48%	0.06%	235
EDDM	0.00%	95.18%	0.00%	15
HDDM	1.65%	76.61%	3.23%	1271
KSWIN	3.32%	42.54%	6.16%	1866
PH	5.18%	72.19%	9.66%	1867
RDDM	0.89%	25.99%	1.72%	2027
STEPD	0.01%	<b>95.80%</b>	0.03%	516

### 5.2. Detailed comparison within each scenario

Tables 6 and 7 present the metrics for Single-Class Local and Global concept drift scenarios. When we analyze the detection performance, it is evident that ADWIN and PH consistently exhibited the best detection performance for both local and global concept drifts, while EDDM, ECDD, and STEPD displayed the least favorable results. Comparing the differences in results between local and global concept drift, we observe that precision values remained relatively stable. However, recall values increased for most drift detectors, with a 17% increase for ADWIN and a 10% increase for PH. This indicates an inverse correlation between false alarms and the magnitude of the concept drift. On the other hand, the detection rate improved from local to global drifts. This is understandable since global concept drifts have a more substantial impact on data distribution, making them easier to detect due to their significant influence.

When considering Multi-Class concept drift scenarios, ADWIN and PH again demonstrated the best results for both scenarios. Similar to the observations in Single-Class concept drift, when we compare local and global drifts, we notice that global drifts were detected more effectively due to their more substantial impact on classifier performance. Furthermore, it is noteworthy that Multi-Class Local drifts exhibited higher recall values than their single-class counterparts. This is expected since a larger portion of the data distribution is influenced by the drift in multi-class scenarios, making it easier to detect and leading to improved recall values.

Additionally, Figs. 6 to 9 illustrate four examples of the data distribution, accuracy of the classifier, and the moments when the ADWIN (the most effective overall) signaled drifts. These visualizations provide insight into

the ease or difficulty in detecting specific scenarios. In the case of Local drifts (Figs. 6 and 8), the drift alerts were signaled when no actual drift occurred, indicating that detecting Local drifts is challenging due to their subtle nature. Conversely, in Global drifts (Figs. 7 and 9), a more distinct alteration in the data distribution was observable, making it easier to identify and leading to more noticeable changes in classifier accuracy. Both Single and Multi scenarios exhibited false alarms, where drift detection alerts were raised despite no actual drift occurring during those periods.

To address **RQ2**, ADWIN and PH consistently demonstrated superior detection performance across all evaluated scenarios. Conversely, EDDM and STEPД, while achieving high recall values, also triggered numerous false alarms, resulting in lower precision. It is worth noting that DDM consistently exhibited low values of both recall and precision in all scenarios, indicating a more conservative approach. Furthermore, when comparing the behavior of drift detectors in the four scenarios, a hierarchy of difficulty emerges, from easiest to hardest detection as follows: Multi-Class Global, Single-Class Global, Multi-Class Local, Single-Class Local. This order of difficulty corresponds closely to the impact on data distribution, with a more pronounced impact leading to a more significant change in error distribution. Another noteworthy observation is that in all evaluated drift detectors, scenarios with more localized drifts tended to generate a higher number of false alarms. This highlights the need for the development of drift detectors that are less sensitive to error distribution, particularly in scenarios with local drifts. In conclusion, our analysis provides insights into how each scenario affected the performance of drift detectors, addressing RQ3.

Table 6: Comparison between drift detectors considering single class local concept drifts.

Drift Detector	Precision	Recall	F1	Delay
ADWIN	<b>5.48%</b>	22.40%	8.80%	2728
DDM	1.88%	5.47%	2.80%	2287
ECDD	0.03%	29.17%	0.07%	201
EDDM	0.00%	<b>97.40%</b>	0.00%	4
HDDM	1.69%	67.45%	3.30%	1359
KSWIN	2.91%	31.51%	5.32%	1991
PH	5.21%	64.06%	<b>9.63%</b>	2176
RDDM	1.13%	31.77%	2.18%	2016
STEPД	0.02%	91.15%	0.04%	649

Table 7: Comparison between drift detectors considering single class global concept drifts.

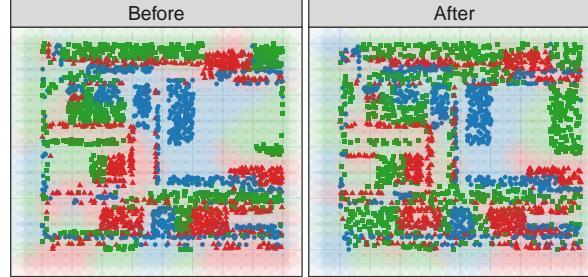
Drift Detector	Precision	Recall	F1	Delay
ADWIN	<b>8.15%</b>	39.06%	<b>13.48%</b>	1909
DDM	2.21%	6.25%	3.27%	2286
ECDD	0.02%	26.30%	0.05%	315
EDDM	0.00%	<b>91.41%</b>	0.00%	11
HDDM	1.81%	69.01%	3.52%	1321
KSWIN	3.30%	34.11%	6.02%	1879
PH	5.87%	75.78%	10.90%	1737
RDDM	1.01%	28.39%	1.96%	1982
STEPD	0.01%	93.75%	0.02%	517

Table 8: Comparison between drift detectors considering multi class local concept drifts.

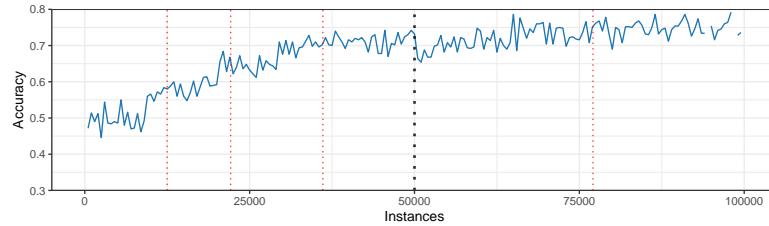
Drift Detector	Precision	Recall	F1	Delay
ADWIN	<b>7.19%</b>	30.40%	<b>11.63%</b>	2360
DDM	1.58%	4.67%	2.36%	2036
ECDD	0.04%	40.29%	0.07%	192
EDDM	0.00%	<b>95.51%</b>	0.00%	21
HDDM	1.67%	80.68%	3.28%	1207
KSWIN	3.31%	44.32%	6.16%	1873
PH	5.18%	70.88%	9.66%	1956
RDDM	0.92%	26.37%	1.77%	2083
STEPD	0.02%	96.70%	0.03%	508

Table 9: Comparison between drift detectors considering multi class global concept drifts.

Drift Detector	Precision	Recall	F1	Delay
ADWIN	<b>7.95%</b>	42.92%	<b>13.41%</b>	2033
DDM	1.94%	6.05%	2.93%	2195
ECDD	0.03%	39.38%	0.06%	277
EDDM	0.00%	<b>95.43%</b>	0.00%	14
HDDM	1.56%	78.88%	3.06%	1299
KSWIN	3.48%	48.86%	6.49%	1819
PH	4.91%	75.80%	9.22%	1706
RDDM	0.72%	21.92%	1.39%	1977
STEPD	0.01%	97.60%	0.03%	469

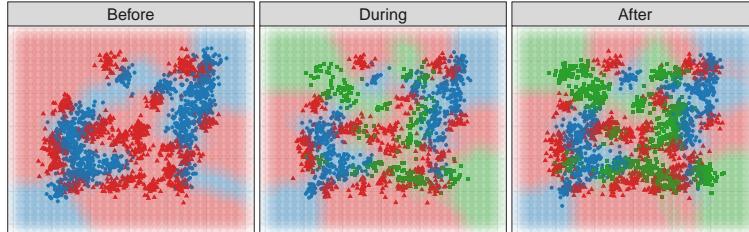


(a) Data distribution before and after a concept drift where new boundaries for the green class are created.

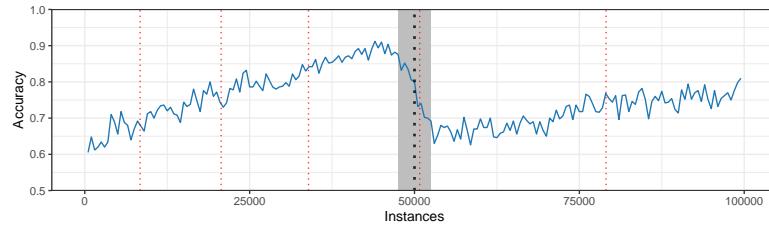


(b) Hoeffing Tree accuracy over time for a specific stream experiencing concept drift. The gray rectangle denotes the duration of the drift. Red dashed lines represent the drift alarms triggered by ADWIN.

Figure 6: Single-Class Local sudden drift data distribution and drift alerts over time considering the `emerging_branch` difficulty and ADWIN as drift detector.

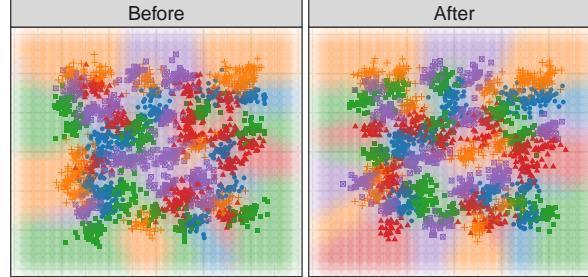


(a) Data distribution before, during and after a concept drift where a new class (green) emerges.

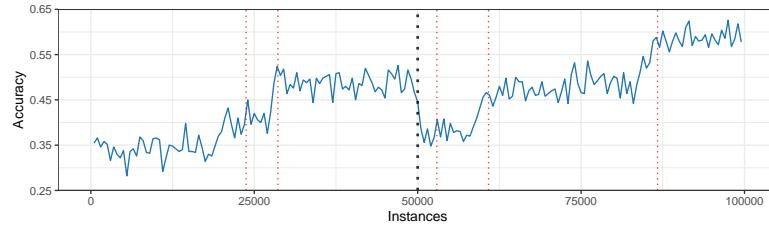


(b) Hoeffing Tree accuracy over time for a specific stream experiencing concept drift. The gray rectangle denotes the duration of the drift. Red dashed lines represent the drift alarms triggered by ADWIN.

Figure 7: Single-Class Global drift data distribution and drift alerts over time considering the `class_emerging` difficulty and ADWIN as drift detector.

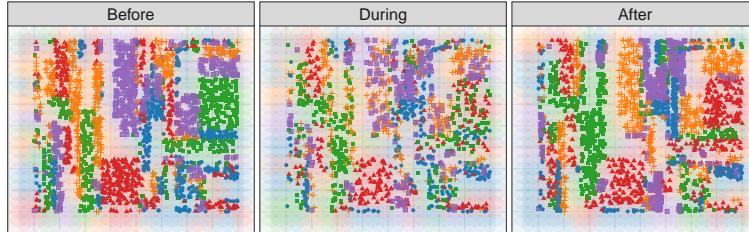


(a) Data distribution before and after a concept drift that swaps positions of some clusters of all classes.



(b) Hoeffing Tree accuracy over time for a specific stream experiencing concept drift. The gray rectangle denotes the duration of the drift. Red dashed lines represent the drift alarms triggered by ADWIN.

Figure 8: Multi-Class Local drift data distribution and drift alerts over time considering the `swap_cluster` difficulty and ADWIN as drift detector.



(a) Data distribution before, during and after a concept drift that swaps positions of all classes.



(b) Hoeffing Tree accuracy over time for a specific stream experiencing concept drift. The gray rectangle denotes the duration of the drift. Red dashed lines represent the drift alarms triggered by ADWIN.

Figure 9: Multi-Class Global drift data distribution and drift alerts over time considering the `swap_leaves` difficulty and ADWIN as drift detector.

### *5.3. Comparison between difficulties*

In this experiment, our aim is to assess the varying levels of difficulty encountered by each drift detector when dealing with different concept drift scenarios. We seek to analyze which scenarios pose more significant challenges and which ones are relatively easier to detect for each drift detector. To accommodate space limitations within the manuscript, we will primarily focus on presenting the results for the top-performing drift detectors, ADWIN and Page Hinkley. The complete set of results is accessible on our website for a more detailed examination.

When considering Single-Class Local drifts, Table 10 provides an overview of the performance metrics for each difficulty as observed with both drift detectors. It is worth noting that, for both detectors, the `reappearing_cluster` scenario appeared to be the easiest to detect. This is because the disappearance and reappearance of data segments tend to induce substantial changes in error distribution, making them relatively clear for the drift detectors, even if they occur in specific regions of the feature space. Conversely, when faced with more complex drift scenarios, ADWIN displayed its worst performance in detecting the `emerging_branch` difficulty, while PH struggled with the `moving_cluster` scenario. The differences in their performance can be attributed to their distinct detection mechanisms. ADWIN employs a two-window approach that relies on detecting a significant difference between the two windows to identify concept drift. Consequently, when the classifier rapidly adapts to the new data distribution, ADWIN may encounter challenges in detecting the drift. On the other hand, PH utilizes a cumulative approach, allowing it to handle such scenarios more effectively. Additionally, difficulties that inherently involve incremental changes, such as those related to moving, merging, and splitting clusters, displayed a higher detection delay in general.

Furthermore, Table 11 provides an overview of the results for each of the Single-Class Global difficulties. Notably, for both ADWIN and PH, the `class_emerging` difficulty appeared to be the least challenging to detect and exhibited the lowest detection delay. This observation aligns with the logic that the emergence of a new class results in a sudden drop in accuracy, as a novel class is being introduced and evaluated. Conversely, the most challenging difficulty to detect was the `moving_cluster`, followed by `splitting_cluster` and `merging_cluster`. These three difficulties are inherently incremental in nature, which means they do not create an abrupt change in the learning curve and error rate. Consequently, this incremental

nature makes them harder to detect promptly. Comparing the results to those of Single-Class Local drifts, an improvement in recall values is noticeable. This suggests that, for almost all difficulties, detecting Global drifts was generally easier due to their broader impact on the data distribution.

Analyzing the metrics for Multi-Class Local drifts in Table 12, ADWIN achieved the best results when encountering the `prune_growth_new_branch` difficulty, while PH exhibited its best performance with `swap_clusters`. Notably, results for difficulties that do not involve incremental drifts were quite similar to the best results. This observation emphasizes the challenge that drift detectors face in effectively detecting incremental drifts, even when multiple classes are affected. Another aspect worth noting is the larger performance gap between ADWIN and PH in the context of multi-class drifts compared to single-class drifts. This difference can be attributed to the higher precision that ADWIN demonstrates in multi-class drift scenarios. When we make comparisons with their Single-Class Local counterparts, it becomes evident that multi-class drifts generally yield better detection results, as expected due to their greater impact on data distribution.

The results for Multi-Class Global drifts are presented in Table 13. ADWIN displayed its best detection performance when confronted with `prune_growth_new_branch` drifts, but it also achieved the satisfactory F1 values for other difficulties, such as `prune_regrowth_branch` and `swap_leaves`. The poorest results for ADWIN once again correlated with difficulties involving incremental drift. PH, on the other hand, exhibited its best performance with `reappearing_cluster`, although its performance with other difficulties was quite similar and generally less favorable than that of ADWIN. In comparison to Multi-Class Local drifts, precision values remained similar, but recall values increased, indicating improved detection even with a notable number of false alarms.

To address **RQ4**, it becomes evident that sudden changes in data distribution, like new data appearing in the feature space or the emergence of a new class, were the easiest and quickest to be detected. Conversely, difficulties that inherently involved incremental changes displayed higher values of False Negatives, resulting in lower recall and less effective detection. It is important to note that this observation is correlated with how the base learner operates. Consequently, when evaluating drift detectors with other base learners or employing unsupervised drift detectors, the ranking of difficulties in terms of complexity may vary.

Table 10: Comparison of detection performance for Single-Class Local drift difficulties by ADWIN and Page Hinkley. Bold values represent the highest value of that column (easiest detection) and underscored values represent the lowest value (harder detection).

Difficulty	Precision		Recall		F1		Delay	
	ADWIN	PH	ADWIN	PH	ADWIN	PH	ADWIN	PH
emerging_branch	<u>3.83%</u>	5.24%	<u>14.58%</u>	<b>70.83%</b>	<u>6.06%</u>	9.76%	2176	1975
emerging_cluster	5.61%	5.41%	25.00%	62.50%	9.16%	9.95%	3364	2249
merging_cluster	5.15%	5.39%	20.83%	62.50%	8.26%	9.92%	2914	2285
moving_cluster	5.03%	<u>4.94%</u>	20.83%	<u>58.33%</u>	8.10%	<u>9.11%</u>	2930	2050
prune_growth_new_branch	4.47%	5.09%	16.67%	<u>66.67%</u>	7.05%	9.45%	2257	2067
prune_regrowth_branch	7.14%	4.94%	22.92%	64.58%	10.89%	9.17%	2479	2261
reappearing_cluster	<b>7.35%</b>	<b>5.56%</b>	<b>37.50%</b>	66.67%	<b>12.29%</b>	<b>10.26%</b>	2550	2179
splitting_cluster	4.95%	5.16%	20.83%	60.42%	8.00%	9.51%	2936	2368

Table 11: Comparison of detection performance for Single-Class Global drift difficulties by ADWIN and Page Hinkley. Bold values represent the highest value of that column (easiest detection) and underscored values represent the lowest value (harder detection).

Difficulty	Precision		Recall		F1		Delay	
	ADWIN	PH	ADWIN	PH	ADWIN	PH	ADWIN	PH
class_emerging_rbf	10.07%	<b>7.68%</b>	<b>58.33%</b>	<b>93.75%</b>	17.18%	<b>14.20%</b>	1104	1698
class_emerging_rt	<b>11.06%</b>	6.55%	52.08%	85.42%	<b>18.25%</b>	12.17%	1198	1280
merging_cluster	5.24%	5.17%	20.83%	60.42%	8.37%	9.52%	2895	2130
moving_cluster	<u>4.37%</u>	<u>4.32%</u>	20.83%	<u>54.17%</u>	<u>7.22%</u>	<u>8.00%</u>	2885	2089
prune_growth_new_branch	9.60%	5.68%	39.58%	77.08%	<u>15.45%</u>	10.59%	2378	1689
prune_regrowth_branch	9.21%	5.54%	43.75%	83.33%	15.22%	10.39%	1722	1487
reappearing_cluster	9.62%	6.76%	58.33%	91.67%	16.52%	12.59%	2124	1896
splitting_cluster	4.50%	5.22%	<u>18.75%</u>	60.42%	7.26%	9.60%	2985	1899

Table 12: Comparison of detection performance for Multi-Class Local drift difficulties by ADWIN and Page Hinkley. Bold values represent the highest value of that column (easiest detection) and underscored values represent the lowest value (harder detection).

Difficulty	Precision		Recall		F1		Delay	
	ADWIN	PH	ADWIN	PH	ADWIN	PH	ADWIN	PH
emerging_branch	<b>11.53%</b>	5.35%	45.37%	76.85%	18.39%	10.01%	2474	1644
emerging_cluster	6.82%	5.15%	27.78%	64.81%	10.95%	9.54%	3016	2329
merging_cluster	4.65%	5.14%	19.44%	64.81%	7.50%	9.52%	2783	2032
moving_cluster	4.66%	4.86%	20.37%	63.89%	7.59%	9.03%	2799	2042
prune_growth_new_branch	11.40%	5.37%	<b>49.07%</b>	80.56%	<b>18.50%</b>	10.08%	1983	1618
prune_regrowth_branch	4.37%	<u>3.88%</u>	15.00%	<u>58.33%</u>	6.77%	<u>7.27%</u>	2650	1970
reappearing_cluster	11.37%	5.69%	48.33%	73.33%	18.41%	10.56%	2696	1749
split_node	<u>3.32%</u>	4.50%	<u>12.04%</u>	61.11%	<u>5.21%</u>	8.38%	1777	2387
splitting_cluster	4.82%	5.58%	20.37%	67.59%	7.80%	10.31%	3013	2213
swap_cluster	7.76%	<b>5.97%</b>	42.59%	<b>85.19%</b>	13.12%	<b>11.15%</b>	1995	1842
swap_leaves	8.23%	5.23%	35.19%	78.70%	13.33%	9.80%	1675	1835

Table 13: Comparison of detection performance for Multi-Class Global drift difficulties by ADWIN and Page Hinkley. Bold values represent the highest value of that column (easiest detection) and underscored values represent the lowest value (harder detection).

Difficulty	Precision		Recall		F1		Delay	
	ADWIN	PH	ADWIN	PH	ADWIN	PH	ADWIN	PH
merging_cluster	<u>4.37%</u>	4.53%	<u>20.37%</u>	60.19%	<u>7.20%</u>	8.43%	2917	1889
moving_cluster	4.93%	<u>4.30%</u>	27.78%	66.67%	8.37%	<u>8.08%</u>	3148	1885
prune_growth_new_branch	<b>11.52%</b>	5.19%	63.89%	87.04%	<b>19.52%</b>	9.80%	1666	1540
prune_regrowth_branch	10.19%	5.14%	55.00%	90.00%	17.19%	9.72%	1865	1485
reappearing_cluster	9.31%	<b>5.51%</b>	58.33%	85.00%	16.06%	<b>10.34%</b>	1797	1705
split_node	7.24%	4.79%	25.93%	67.59%	11.31%	8.94%	2129	2061
splitting_cluster	4.97%	4.62%	20.37%	<u>58.33%</u>	7.99%	8.55%	2987	2177
swap_cluster	8.76%	5.12%	<b>69.44%</b>	<b>91.67%</b>	15.56%	9.70%	1847	1263
swap_leaves	9.79%	5.13%	57.41%	86.11%	16.73%	9.68%	1657	1608

#### 5.4. Impact of number of classes and number of features

Furthermore, the benchmark dataset used in the experiments included streams with varying configurations in terms of the number of classes and features. Therefore, we assessed how these different specifications influenced drift detection while considering each drift categorization to address **RQ5**. Fig. 10 presents a boxplot showing the F1 score for ADWIN and Page Hinkley in each drift category, considering different numbers of classes, and Fig. 11 considering different feature space dimensionality.

Firstly, it is evident that as the number of classes in the stream increases, the detection of concept drifts becomes more challenging, regardless of how the stream is affected. Moreover, when considering Single-Class Global drifts, a noticeable drop in performance is observed. This is expected in this specific scenario because when only one class is affected, the more stationary the other classes are, the more localized the change becomes, making it difficult to detect. On the other hand, when multiple classes are affected globally, the increase in the number of classes did not have a significant impact on drift detection. The scenario that posed the most challenge for both drift detectors was Single-Class Local drift with only 2 classes. In a binary setting where only a portion of the classes is affected, it is similar to an imbalanced scenario, and these changes may not significantly impact the classifier’s error.

When considering different feature space configurations, we can see that this factor did not have a substantial impact on Page Hinkley’s performance. However, it did affect ADWIN, particularly in the presence of Single-Class Local and Multi-Class Local drifts. Additionally, it is worth noting that

ADWIN's performance exhibits higher variance, while Page Hinkley demonstrates consistent performance across various stream configurations.

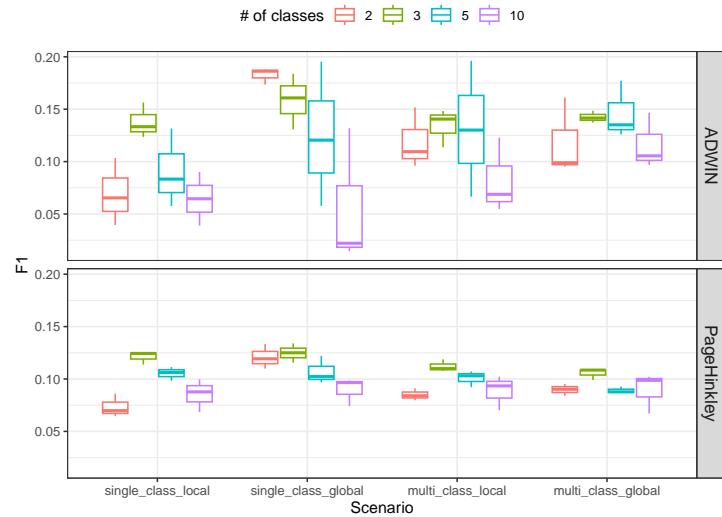


Figure 10: F1 Score for each categorization of concept drift. Each color represents streams with different number of classes.

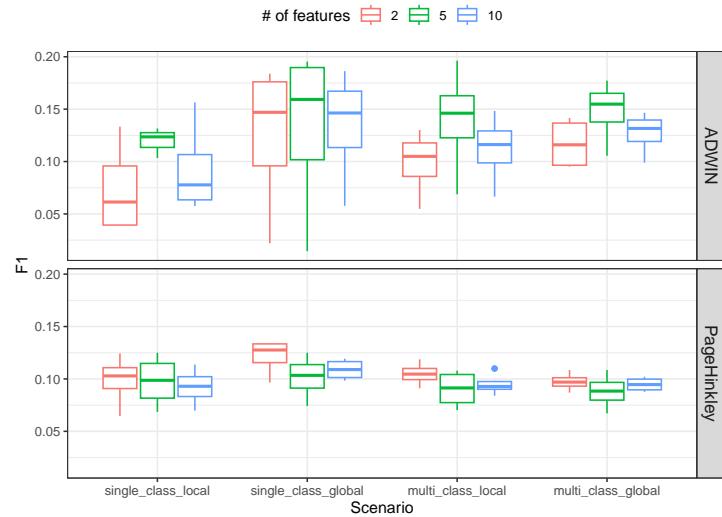


Figure 11: F1 Score for each categorization of concept drift. Each color represents streams with different number of features.

### *5.5. Impact on classifier performance*

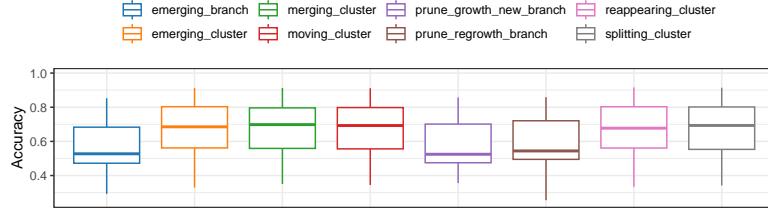
Finally, in addition to assessing how the proposed scenarios and difficulties influence concept drift detection, we also examined their impact on the performance of the utilized classifier, the Hoeffding Tree, when used in conjunction with the best-performing drift detector. The primary objective of this experiment is to ascertain how each of the proposed scenarios affects classification performance and to determine whether combining them with a drift detector enhances overall performance.

Figs. 12 and 13 present the accuracy corresponding to each drift difficulty. Analyzing the average accuracy enables identification of the drifts causing the most significant drop in accuracy, given that, prior to the drift, all data streams from the same generator are equal.

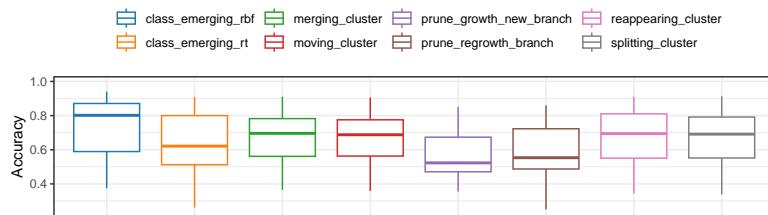
In the Single-Class scenario, particularly focusing on local drifts, incremental changes such as emerging, merging, moving, and splitting clusters proved to be the most challenging to detect. However, classifiers exhibited a capacity for self-adaptation and showed their best accuracy levels. The emergence of new branches posed the most complex scenario in terms of predictive performance. Furthermore, streams generated with Random Tree generator displayed lower accuracy compared to those with RBF. In terms of global changes, the results were similar to local changes, except for `class_emerging`, which exhibited higher accuracy values. This is likely due to the reduction in the number of classes for the majority of the stream, leading to improved accuracy.

Within the Multi-Class local and global scenarios, performance aligned closely with what was observed in Single-Class drifts. Incremental drifts demonstrated the highest accuracy levels, whereas `emerging_branch` and `swap_leaves` presented the most challenging difficulties. Once again, the Random Tree generator proved more challenging to learn from than RBF, which could be attributed to the distribution of clusters. As feature space dimensions increased, the RBF model facilitated easier learning of class boundaries.

Furthermore, we conducted a comparison to understand the impact of integrating the best-performing drift detector, i.e. ADWIN, on the classifier's performance. For this analysis, we compared the standard classifier with two modified versions: Drift/Warning Retraining (HT-DW) and Adaptive Hoeffding Tree (AHT). The former initiates retraining of the classifier upon the detector signaling a warning and completely replaces the classifier when a drift is detected. In contrast, the latter employs a drift detector for each

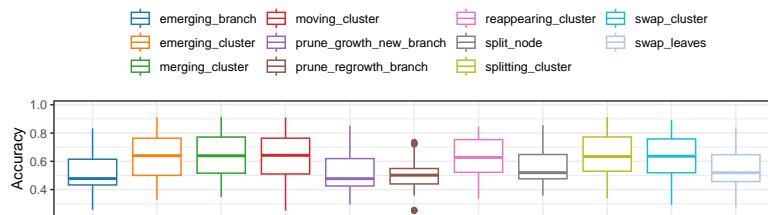


(a) Single-Class Local

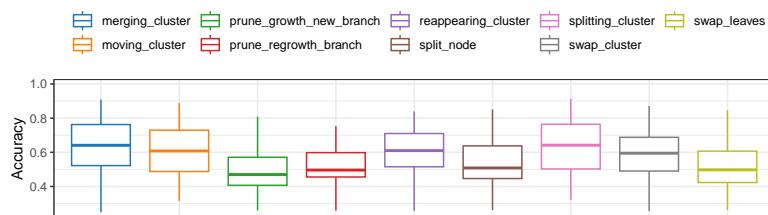


(b) Single-Class Global

Figure 12: Boxplot of HT accuracy for each of the proposed difficulties for Single-Class Local (a) and Global (b)



(a) Multi-Class Local



(b) Multi-Class Global

Figure 13: Boxplot of HT accuracy for each of the proposed difficulties for Multi-Class Local (top) and Global (bottom)

Table 14: Accuracy of HT and AHT for evaluated drift categories.

Drift Category	HT	AHT	HT-DW
Single-Class Local	66.88% $\pm$ 0.148	70.71% $\pm$ 0.127 ( $\uparrow$ 3.83%)	50.13% $\pm$ 0.135 ( $\downarrow$ 16.75%)
Single-Class Global	68.16% $\pm$ 0.148	71.61% $\pm$ 0.132 ( $\uparrow$ 3.44%)	52.55% $\pm$ 0.157 ( $\downarrow$ 15.61%)
Multi-Class Local	60.78% $\pm$ 0.147	65.18% $\pm$ 0.126 ( $\uparrow$ 4.39%)	44.92% $\pm$ 0.121 ( $\downarrow$ 15.87%)
Multi-Class Global	58.80% $\pm$ 0.144	63.28% $\pm$ 0.128 ( $\uparrow$ 4.48%)	45.52% $\pm$ 0.122 ( $\downarrow$ 13.28%)

branch of the tree, replacing only the affected branch when a drift signal is detected, while maintaining the overall structure of the tree. This differentiates HT-DW’s global approach from AHT’s localized strategy in addressing concept drift within the classifier structure.

Table 14 displays the average accuracy for each drift category across the three classifiers. Overall, the AHT strategy, focusing on addressing drift locally, improved the accuracy of the standard Hoeffding Tree by an average of 4.03%. This improvement underscores the robustness of the Hoeffding Tree in adapting to new concepts, exhibiting good accuracy values even without a drift detector. On the other hand, despite the utilization of a drift detector, completely retraining the classifier resulted in decreased accuracy across all evaluated scenarios. Notably, the more localized the drift category, the more pronounced the reduction in accuracy.

Additionally, in Figs. 14-17, the accuracy of each classifier under the `prune_growth_new_branch`, `reappearing_cluster`, `reappearing_cluster` and `reappearing_cluster` difficulties is displayed, including sudden and gradual drifts. These figures serve as an illustration of the classifiers’ behavior across different drift categories. Across Single-Class drifts, both local and global, AHT consistently demonstrated superior accuracy. Even in the presence of drift, HT showcased a remarkable self-adaptation ability, recovering swiftly. Conversely, the retraining strategy resulted in notably poorer performance. Notably, the sole scenario where retraining displayed an advantage was in the context of Multi-Class Global drifts. This specific category provoked a significant change in data distribution, evident from the notable accuracy drop. Given the entirely new data distribution, retraining the classifier proved a valid method for recovering from these drifts. However, by the end of the stream, the plain HT exhibited better accuracy compared to both drift-adapted classifiers.

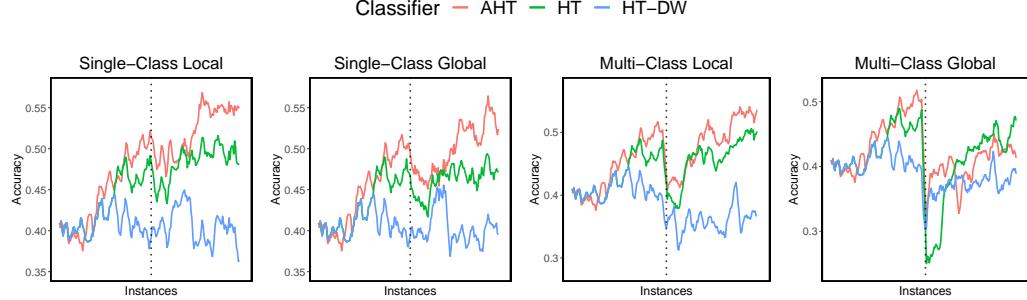


Figure 14: Accuracy curve for AHT, HT and HT-DW on `prune_growth_new_branch` data stream in the presence of a sudden concept drift. Dashed vertical line indicates where the concept drift happened.

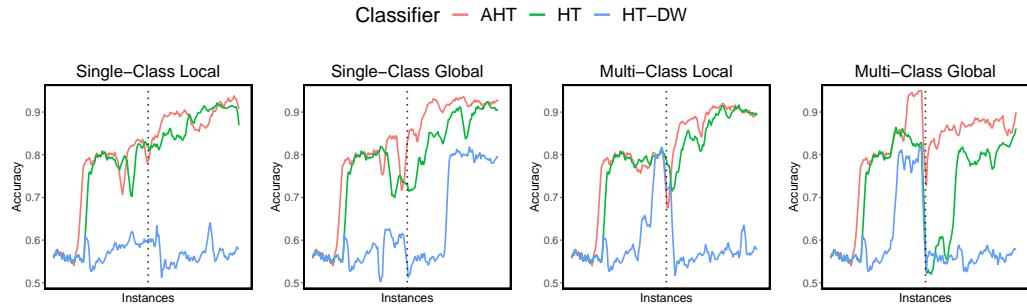


Figure 15: Accuracy curve for AHT, HT and HT-DW on `reappearing_cluster` data stream in the presence of a sudden concept drift. Dashed vertical line indicates where the concept drift happened.

## 6. Lessons learned and recommendations

In order to summarize the knowledge we extracted through the experimental evaluation, this section presents the lessons learned and recommendations for future researchers.

**Locality matters.** The locality of concept drift plays a pivotal role in the accuracy and adaptability of detectors when detecting changes in data streams. When dealing with the detection of concept drift, understanding its locality provides insights into how changes within the data stream, affecting either the entire stream or specific parts of it, influence the classification system. For instance, in scenarios where drift affects only particular sections, such as individual clusters or subsets of classes, it is critical to discern the

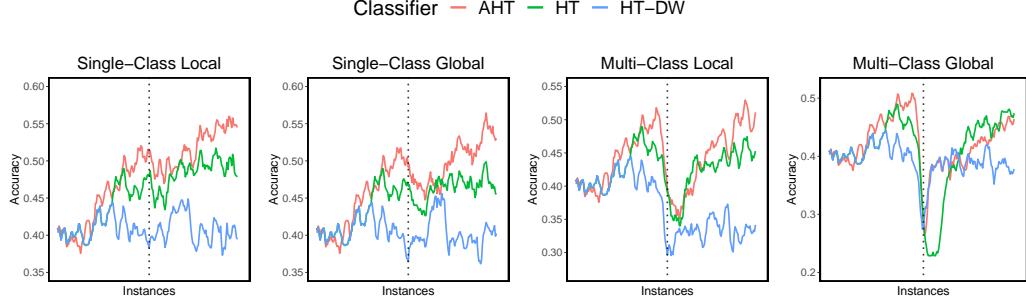


Figure 16: Accuracy curve for AHT, HT and HT-DW on `prune_growth_new_branch` data stream in the presence of a gradual concept drift. Dashed vertical line indicates where the concept drift happened.

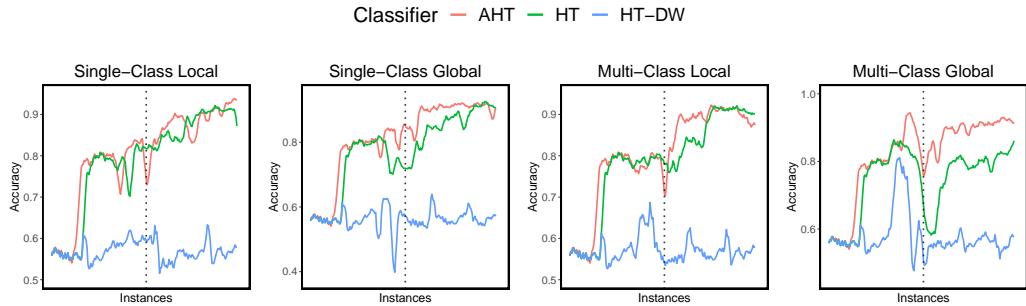


Figure 17: Accuracy curve for AHT, HT and HT-DW on `reappearing_cluster` data stream in the presence of a gradual concept drift. Dashed vertical line indicates where the concept drift happened.

precise areas within the data stream that undergo these changes. Different sections might require varied adaptations to maintain predictive performance. Moreover, different types of drift, local or global, might require different strategies to handle them effectively. Understanding these subtleties can significantly impact the accuracy and performance of the drift detection and adaptation methods used in learning models.

**Handling high number of classes.** The learning task becomes more challenging for both classifiers and drift detectors as the number of classes increases. Our study indicated that when the impact is localized due to fewer classes being affected, the detection of the change becomes more difficult. Hence, drift detectors capable of handling a larger number of classes

play a crucial role in effective drift detection. Furthermore, future research should aim at detecting drifts within individual classes, a strategy that could facilitate the adaptation of classifiers to new concepts.

**False alarms.** Drift detectors that rely on error rates often generate numerous false alarms. Given their univariate nature, even positive changes in the error rate, such as when a new branch is created in an incremental tree or a classifier establishes a new boundary, can trigger alarms in certain drift detectors, which is counterproductive. In more complex scenarios, such as localized changes, false alarms can mislead the classifier, resulting in poorer performance. For future research, it is essential that drift detectors strike a better balance between detection and false alarms by examining not only the error curve but also statistics related to data distribution.

The common strategy to address concept drift involves replacing the classifier once a drift is identified. While this approach works well with ensembles due to their varied classifiers, our experiments have shown that when handling drifts that only affect specific boundaries or when the classifier is capable of self-adaptation, avoiding complete replacement of the established boundaries results in better performance. Furthermore, retraining the classifier for each identified drift may lead to redundant work due to the high occurrence of false alarms, even with state-of-the-art drift detectors. This emphasizes the necessity for drift detectors capable of pinpointing which class or segment of the data stream was affected by the drift. This, in turn, facilitates a more precise adaptation of the classifier rather than a complete replacement.

**Retraining is not (always) the best option.** The common strategy to address concept drift involves replacing the classifier once a drift is identified. While this approach works well with ensembles due to their varied classifiers, our experiments have shown that when handling drifts that only affect specific boundaries or when the classifier is capable of self-adaptation, avoiding complete replacement of the established boundaries results in better performance. Furthermore, retraining the classifier for each identified drift may lead to redundant work due to the high occurrence of false alarms, even with state-of-the-art drift detectors. This emphasizes the necessity for drift detectors capable of pinpointing which class or segment of the data stream was affected by the drift. This, in turn, facilitates a more precise adaptation of the classifier rather than a complete replacement.

## 7. Conclusion and future work

In this paper, we presented a comprehensive study focusing on benchmarking and evaluating the impact of concept drift, specifically concerning its locality and magnitude, on classifiers and drift detectors. We introduced a novel categorization of concept drift, considering its locality and scale. Through a systematic approach, we identified significant challenges within this domain and formulated a set of 2,760 benchmark problems that encompass various levels of difficulty, guided by our proposed categorization. Additionally, we conducted a comparative evaluation of 9 state-of-the-art drift detectors across a diverse range of difficulties. Each drift detector was thoroughly assessed across various isolated and combined scenarios. This analysis not only identifies the top-performing detectors but also sheds light on their specific strengths, providing valuable insights for future research in drift detection. Moreover, we evaluated how the locality of drift influences the base-classifier’s performance, gaining knowledge into the most effective approaches for addressing different categories of concept drift to minimize recovery time. This comprehensive study offered insights into different and yet unexplored categorizations of concept drift, providing an understanding of how drift detectors and classifiers perform across numerous difficulties. All benchmark problems and evaluation methodologies are publicly available, enabling future researchers to design efficient approaches, drift detectors, and classifiers to handle local concept drifts effectively.

Our future work aims to explore the influence of concept drift locality on data streams characterized by imbalanced class distributions. Additionally, our plans include an investigation into the performance of unsupervised and semi-supervised drift detectors within the proposed scenarios. Furthermore, we seek to develop classifiers with the capacity to handle concept drift locally, avoiding the necessity for retraining or complete replacement.

## References

- [1] J. Gama, P. Medas, G. Castillo, P. Rodrigues, Learning with drift detection, in: Brazilian Symposium on Artificial Intelligence, 2004.
- [2] M. Bahri, A. Bifet, J. Gama, H. M. Gomes, S. Maniu, Data stream analysis: Foundations, major tasks and tools, Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery (2021).

- [3] J. Gama, Knowledge discovery from data streams, CRC Press, 2010.
- [4] J. Gama, I. Žliobaitė, A. Bifet, M. Pechenizkiy, A. Bouchachia, A survey on concept drift adaptation, ACM Computing Surveys (2014).
- [5] G. Aguiar, B. Krawczyk, A. Cano, A survey on learning from imbalanced data streams: taxonomy, challenges, empirical study, and reproducible experimental framework, Machine Learning (2023).
- [6] L. Korycki, B. Krawczyk, Concept drift detection from multi-class imbalanced data streams, in: IEEE International Conference on Data Engineering (ICDE), 2021.
- [7] A. D. Viniski, J. P. Barddal, A. de Souza Britto Jr, F. Enembreck, H. V. A. de Campos, A case study of batch and incremental recommender systems in supermarket data under concept drifts and cold start, Expert Systems with Applications (2021).
- [8] A. L. Suárez-Cetrulo, D. Quintana, A. Cervantes, A survey on machine learning for recurring concept drifting data streams, Expert Systems with Applications (2023).
- [9] N. Lu, J. Lu, G. Zhang, R. L. De Mantaras, A concept drift-tolerant case-base editing technique, Artificial Intelligence (2016).
- [10] A. Liu, Y. Song, G. Zhang, J. Lu, Regional concept drift detection and density synchronized drift adaptation, in: International Joint Conference on Artificial Intelligence (IJCAI), 2017.
- [11] R. S. Barros, D. R. Cabral, P. M. Gonçalves Jr, S. G. Santos, Rddm: Reactive drift detection method, Expert Systems with Applications (2017).
- [12] E. B. Gulcan, F. Can, Unsupervised concept drift detection for multi-label data streams, Artificial Intelligence Review (2023).
- [13] J. Gama, G. Castillo, Learning with local drift detection, in: International conference on advanced data mining and applications, Springer, 2006.

- [14] B. Krawczyk, L. L. Minku, J. Gama, J. Stefanowski, M. Woźniak, Ensemble learning for data stream analysis: A survey, *Information Fusion* (2017).
- [15] J. Lu, A. Liu, F. Dong, F. Gu, J. Gama, G. Zhang, Learning under concept drift: A review, *IEEE Transactions on Knowledge and Data Engineering* (2018).
- [16] A. R. Masegosa, A. M. Martínez, D. Ramos-López, H. Langseth, T. D. Nielsen, A. Salmerón, Analyzing concept drift: A case study in the financial sector, *Intelligent Data Analysis* (2020).
- [17] G. I. Webb, R. Hyde, H. Cao, H. L. Nguyen, F. Petitjean, Characterizing concept drift, *Data Mining and Knowledge Discovery* (2016).
- [18] B. Krawczyk, A. Cano, Online ensemble learning with abstaining classifiers for drifting and noisy data streams, *Applied Soft Computing* (2018).
- [19] G. Ditzler, M. Roveri, C. Alippi, R. Polikar, Learning in nonstationary environments: A survey, *IEEE Computational Intelligence Magazine* (2015).
- [20] R. S. M. Barros, S. G. T. C. Santos, A large-scale comparison of concept drift detectors, *Information Sciences* (2018).
- [21] E. S. Page, Continuous inspection schemes, *Biometrika* (1954).
- [22] S. Roberts, Control chart tests based on geometric moving averages, *Technometrics* (2000).
- [23] M. Baena-García, J. del Campo-Ávila, R. Fidalgo, A. Bifet, R. Gavalda, R. Morales-Bueno, Early drift detection method, in: *International Workshop on Knowledge Discovery from Data Streams*, 2006.
- [24] I. Frias-Blanco, J. del Campo-Ávila, G. Ramos-Jimenez, R. Morales-Bueno, A. Ortiz-Díaz, Y. Caballero-Mota, Online and non-parametric drift detection methods based on hoeffding's bounds, *IEEE Transactions on Knowledge and Data Engineering* (2014).
- [25] R. S. M. de Barros, J. I. G. Hidalgo, D. R. de Lima Cabral, Wilcoxon rank sum test drift detector, *Neurocomputing* (2018).

- [26] G. J. Ross, N. M. Adams, D. K. Tasoulis, D. J. Hand, Exponentially weighted moving average charts for detecting concept drift, *Pattern Recognition Letters* (2012).
- [27] A. Bifet, R. Gavalda, Learning from time-changing data with adaptive windowing, in: *SIAM International Conference on Data Mining*, 2007.
- [28] C. Raab, M. Heusinger, F.-M. Schleif, Reactive soft prototype computing for concept drift streams, *Neurocomputing* (2020).
- [29] K. Nishida, K. Yamauchi, Detecting concept drift using statistical testing, in: *International Conference on Discovery Science*, Springer, 2007.
- [30] D. T. J. Huang, Y. S. Koh, G. Dobbie, R. Pears, Detecting volatility shift in data streams, in: *IEEE International Conference on Data Mining*, 2014.
- [31] A. Pesaranghader, H. Viktor, Fast hoeffding drift detection method for evolving data streams, in: *Machine Learning and Knowledge Discovery in Databases: European Conference*, 2016.
- [32] J. Komorniczak, P. Ksieniewicz, Complexity-based drift detection for nonstationary data streams, *Neurocomputing* (2023).
- [33] P. Wang, H. Yu, N. Jin, D. Davies, W. L. Woo, Quadcdd: A quadruple-based approach for understanding concept drift in data streams, *Expert Systems with Applications* (2024).
- [34] A. Lapiński, B. Krawczyk, P. Ksieniewicz, M. Woźniak, An empirical insight into concept drift detectors ensemble strategies, in: *IEEE Congress on Evolutionary Computation*, 2018.
- [35] P. Sobolewski, M. Woźniak, Comparable study of statistical tests for virtual concept drift detection, in: *International Conference on Computer Recognition Systems*, 2013.
- [36] X. Song, M. Wu, C. Jermaine, S. Ranka, Statistical change detection for multi-dimensional data, in: *ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 2007.

- [37] A. A. Qahtan, B. Alharbi, S. Wang, X. Zhang, A pca-based change detection framework for multidimensional data streams: Change detection in multidimensional data streams, in: ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, 2015.
- [38] F. Gu, G. Zhang, J. Lu, C.-T. Lin, Concept drift detection based on equal density estimation, in: International Joint Conference on Neural Networks (IJCNN), IEEE, 2016.
- [39] L. Bu, C. Alippi, D. Zhao, A pdf-free change detection test based on density difference estimation, IEEE Transactions on Neural Networks and Learning Systems (2016).
- [40] A. Liu, J. Lu, F. Liu, G. Zhang, Accumulating regional density dissimilarity for concept drift detection in data streams, Pattern Recognition (2018).
- [41] N. Lu, G. Zhang, J. Lu, Concept drift detection via competence models, Artificial Intelligence (2014).
- [42] P. M. Gonçalves Jr, S. G. de Carvalho Santos, R. S. Barros, D. C. Vieira, A comparative study on concept drift detectors, Expert Systems with Applications (2014).
- [43] S. G. Santos, R. S. Barros, P. M. Gonçalves Jr, A differential evolution based method for tuning concept drift detectors in data streams, Information Sciences (2019).
- [44] E. S. Babüroğlu, A. Durmuşoğlu, T. Dereli, Novel hybrid pair recommendations based on a large-scale comparative study of concept drift detection, Expert Systems with Applications (2021).
- [45] L. Poenaru-Olaru, L. Cruz, A. van Deursen, J. S. Rellermeyer, Are concept drift detectors reliable alarming systems?-a comparative study, in: IEEE International Conference on Big Data, IEEE, 2022.
- [46] G. Y. Sakurai, J. F. Lopes, B. B. Zarpelão, S. Barbon Junior, Benchmarking change detector algorithms from different concept drift perspectives, Future Internet (2023).

- [47] D. Brzezinski, L. L. Minku, T. Pewinski, J. Stefanowski, A. Szumaczuk, The impact of data difficulty factors on classification of imbalanced and concept drifting data streams, Knowledge and Information Systems (2021).
- [48] M. Lango, J. Stefanowski, What makes multi-class imbalanced problems difficult? an experimental study, Expert Systems with Applications (2022).
- [49] G. Holmes, R. Kirkby, B. Pfahringer, Stress-testing hoeffding trees, in: Knowledge Discovery in Databases: PKDD 2005: 9th European Conference on Principles and Practice of Knowledge Discovery in Databases, Springer, 2005.
- [50] J. Montiel, M. Halford, S. M. Mastelini, G. Bolmier, R. Sourty, R. Vaysse, A. Zouitine, H. M. Gomes, J. Read, T. Abdessalem, et al., River: machine learning for streaming data in python, The Journal of Machine Learning Research (2021).