



CURIE: a cellular automaton for concept drift detection

Jesus L. Lobo¹ · Javier Del Ser² · Eneko Osaba¹ · Albert Bifet³ · Francisco Herrera⁴

Received: 19 September 2020 / Accepted: 18 June 2021 / Published online: 4 September 2021

© The Author(s), under exclusive licence to Springer Science+Business Media LLC, part of Springer Nature 2021

Abstract

Data stream mining extracts information from large quantities of data flowing fast and continuously (data streams). They are usually affected by changes in the data distribution, giving rise to a phenomenon referred to as *concept drift*. Thus, learning models must detect and adapt to such changes, so as to exhibit a good predictive performance after a drift has occurred. In this regard, the development of effective drift detection algorithms becomes a key factor in data stream mining. In this work we propose *CURIE*, a drift detector relying on cellular automata. Specifically, in *CURIE* the distribution of the data stream is represented in the grid of a cellular automata, whose neighborhood rule can then be utilized to detect possible distribution changes over the stream. Computer simulations are presented and discussed to show that *CURIE*, when hybridized with other base learners, renders a competitive behavior in terms of detection metrics and classification accuracy. *CURIE* is compared with well-established drift detectors over synthetic datasets with varying drift characteristics.

Keywords Concept drift · Drift detection · Data stream mining · Cellular automata

1 Introduction

Data Stream Mining (DSM) techniques are focused on extracting patterns from continuous (potentially infinite) and fast data. A data stream is the basis of machine learning techniques for this particular kind of data, which is composed of an ordered sequence of instances that arrive one by one or in batches. Depending on the constraints imposed by the application scenario at hand, such instances can be read only once or at most a reduced number of times, using limited computing and memory resources. These constraints require an incremental learning (or one-pass learning) procedure where

Dedicated to Tom Fawcett and J. H. Conway, who passed away in 2020, for their noted contributions to the field of cellular automata and machine learning, and for inspiring this research work.

Responsible editor: Annalisa Appice, Sergio Escalera, Jose A. Gamez, Heike Trautmann.

Extended author information available on the last page of the article

past data cannot be stored for batch training in future time steps. Due to these challenging conditions under which learning must be done, DSM has acquired a notable relevance in recent years, mostly propelled by the advent of Big Data technologies and data-intensive practical use cases (Bifet et al. 2018).

In this context, data streams are often generated by non-stationary phenomena, which may provoke a change in the distribution of the data instances (and/or their annotation). This phenomenon is often referred to as *concept drift* (Webb et al. 2016). These changes cause that predictive models trained over data streams become eventually obsolete, not adapting suitably to the new distribution (concept). The complexity of overcoming this issue, and its prevalence over many real scenarios, make *concept drift* detection and adaptation acknowledged challenges in DSM (Jie et al. 2018). Examples of data stream sources undergoing *concept drift* include computer network traffic, wireless sensor data, phone conversations, social media, marketing data, ATM transactions, web searches, and electricity consumption traces, among others (Žliobaitė 2016). Recently, several emerging paradigms such as the so-called Smart Dust (Ilyas and Mahgoub 2018), Utility Fog (Dastjerdi and Buyya 2016), Microelectromechanical Systems (MEMS or “motes”) (Judy 2001), or Swarm Intelligence and Robotics (Del Ser et al. 2019), are in need for efficient and scalable solutions in real-time scenarios. Here *concept drift* may be present, and thus making drift detection a necessity.

This complexity in the *concept drift* phenomenon manifests when researchers try to characterize it (Webb et al. 2016). Indeed, there are many different types of concept drifts, characterized by e.g. the speed or severity of change. Consequently, drift detection is a key factor for those active strategies that require triggering mechanisms for drift adaptation (Hu et al. 2019). A drift detector estimates the time instants at which changes occur over the stream, so that when a change is detected, an adaptation mechanism is applied to the base learner so as to avoid a degradation of its predictive performance. The design of a *concept drift* detector with high performance is not trivial, yet it is crucial to achieve more reliable DSM models. In fact, a general-purpose strategy for *concept drift* detection, handling and recovery still remains as an open research avenue, as foretold by the fulfillment of the No Free Lunch theorem in this field (Hu et al. 2019). This difficulty to achieve a universal best approach becomes evident in the most recent comparative among drift detectors made in Barros and Santos (2018). Analyzing its mean rank of methods, we observe how there is not a method with the best metrics, or even showing the best performance in the majority of them. In this regard, the design objective is to develop techniques that detect all existing drifts in the stream with low latency and as few false alarms and missed detections as possible. Thus, the most suitable drift detector depends on the characteristics of the DSM problem under study, giving more emphasis to some metrics than others. Regarding the detection metrics, we usually tend to put in value those drift detectors that are able to show a good classification performance while minimizing the distance of the true positive detections.

Cellular automata (CA), as low-bias and robust-to-noise pattern recognition methods with competitive classification performance, meet the requirements imposed by the aforementioned paradigms mainly due to their simplicity and parallel nature. In this work we present a Cellular aUtomaton for concept dRift dEtection (*CURIE*), capable of competitively identifying drifts over data streams. The proposed approach is based

on CA, which became popular when Conway's Game of Life appeared in 1970, and thereafter attracted attention when Stephen Wolfram published his CA study in 2002 (Wolfram 2002). Although CA are not very popular in the data mining field, Fawcett showed in Fawcett (2008) that they can become simple, low-bias methods. *CURIE*, as any other CA-based technique, is computationally complete (able to perform any computation which can be done by digital computers) and can model complex systems from simple structures, which puts it in value to be considered in the DSM field. Moreover, *CURIE* is tractable, transparent and interpretable (Lobo et al. 2021), all ingredients that have lately attracted attention under the eXplainable Artificial Intelligence (XAI) paradigm (Arrieta et al. 2020), and not easy to achieve when designing new data mining techniques. The natural concordance between data and the internal structure of a cellular automaton makes *CURIE* to be closer to a transparent model by design, leaving aside any need for external algorithmic components (post-hoc explainability tools) to interpret its decisions (Rudin 2019). Next, we summarize the main contributions of *CURIE* to the drift detection field:

- It is capable of competitively detecting abrupt and gradual concept drifts.
- It does not require the output (class prediction) of the base learner. Instead, it extracts the required information for drift detection from its internal structure, looking at the changes occurring in the neighborhood of cells.
- It is transparent by design due to the fact that its cellular structure is a direct representation of the feature space and the labels to be predicted.
- It can be combined with any base learner.

Besides, *CURIE* offers another additional advantage in DSM:

- It is also able to act as an incremental learner and adapt to the change (Lobo et al. 2021), going one step further by embodying an *all-in-one* approach (learner and detector).

The rest of the manuscript is organized as follows: first, we provide the background of the field in Sect. 2. Next, we introduce the fundamentals of CA and their role in DSM in Sect. 3. Section 4 exposes the details of our proposed drift detector *CURIE*, whereas Sects. 5 and 6 elaborate on experimental setup and analyze results with synthetic and real-world data stream respectively. Finally, Sect. 7 concludes the manuscript with an outlook towards future research derived from this work.

2 Related work

We now delve into the background literature related to the main topics of this work: drift detection (Sect. 2.1) and cellular automata for machine learning (Sect. 2.2).

2.1 Drift detection

DSM has attracted much attention from the machine learning community (Gomes et al. 2019). Researchers are now on the verge of moving out DSM methods from laboratory environments to real scenarios and applications, similarly to what occurred

with traditional machine learning methods in the past. Most efforts in DSM have been focused on supervised learning (Bifet et al. 2018) (mainly on classification), addressing the *concept drift* problem (Webb et al. 2016). Generally, these efforts have been invested in the development of new methods and algorithms that maintain an accurate decision model with the capability of learning incrementally from data streams while forgetting concepts (Losing et al. 2018).

For this purpose, drift detection and adaptation mechanisms are needed (Jie et al. 2018). In contrast to passive (blind) approaches where the model is continuously updated every time new data instances are received (i.e., drift detection is not required), active strategies (where the model gets updated only when a drift is detected) are in need for effective drift detection mechanisms. Most active approaches usually utilize a specific classifier (base learner) and analyze its classification performance (e.g. accuracy or error rate) to indicate whether a drift has occurred or not. Then, the base learner is trained on the current instance within an incremental learning process repeated for each incoming instance of the data stream. Despite the most used input for the drift detectors are the accuracy or error rate, we can find other detectors that use other inputs such as diversity (Minku and Yao 2011) or structural changes stemming from the model itself (Lobo et al. 2018).

There is a large number of drift detectors in the literature, many of them compared in Gonçalves Jr et al. (2014). As previously mentioned, the conclusion of these and other works is that there is no a general-purpose strategy for *concept drift*. The selection of a good strategy depends on the type of drift and particularities of each data streaming scenario. Other more recent *concept drift* detection mechanisms have been presented and well described in Barros and Santos (2018).

2.2 Cellular automata for pattern recognition

CA are not popular in the pattern recognition community, but even so we can find recent studies and applications. In Collados-Lara et al. (2019), authors propose CA to simulate potential future impacts of climate change on snow covered areas, whereas in Gounaridis et al. (2019) an approach to explore future land use/cover change under different socio-economic realities and scales is presented. Scheduling is another field where CA has been profusely in use (Carvalho and Carneiro 2019). Another recent CA approach for classification is Uzun et al. (2018). CA have been also used with convolutional neural networks (Gilpin 2019) and reservoir computing (Nichele and Molund 2017).

Regarding DSM or *concept drift* detection fields, the presence of CA-based proposals is even scarcer. Although a series of contributions unveiled solid foundations for CA to be used for pattern recognition (Raghavan 1993), it was not until 2008 (Fawcett 2008) [departing from the seminal work in Ultsch (2002)] when CA was presented as a simple but competitive method for parallelism, with a low-bias, effective and competitive in terms of classification performance, and robust to noise. Regarding DSM and *concept drift* detection, timid efforts have been reported so far in Hashemi et al. (2007) and Pourkashani and Kangavari (2008), which must be considered as early attempts to deal with noise rather than with incremental learning and drift detection. They used

a CA-based approach as a real-time instance selector to avoid noisy instances, while the classification task was performed in batch learning mode by non-CA-based learning algorithms. Thus, CA is proposed as a mere complement to select instances, and not as an incremental learner. Besides, their detection approach is simply based on the local class disagreements between neighboring cells, without considering relevant aspects such as the grid size, the radius of the neighborhood, or the moment of the disagreement, among other factors. Above all, they do not provide any evidence on how their solution learns incrementally, nor details on the real operation of the drift detection approach. Finally, in terms of drift detection evaluation, their approach is not compared to known detectors using reputed base learners and standard detection metrics.

More recently, the authors of Lobo et al. (2021) transform a cellular automaton into a real incremental learner with drift adaptation capabilities. In this work, we go one step further by proposing *CURIE*, a cellular automaton featuring a set of novel ingredients that endow it with abilities for drift detection in DSM. As we will present in detail, *CURIE* is an interpretable CA-based drift detector, able to detect abrupt and gradual drifts, and providing very competitive classification performances and detection metrics.

3 Cellular automata

3.1 Foundations

Von Neumann described CA as discrete dynamical systems with a capacity of universal computability (Von Neumann and Burks 1966). Their simple local interaction and computation of cells result in a huge complex behavior when these cells act together, being able to describe complex systems in several scientific disciplines.

Following the notation of Kari in Kari (2005), a cellular automaton can be formally defined as: $A \doteq (d, \mathcal{S}, f_{\boxplus}, f_{\circ})$, with d denoting the dimension, \mathcal{S} a group of discrete states, $f_{\boxplus}(\cdot)$ a function that receives as input the coordinates of the cell and returns the neighbors of the cell to be utilized by the update rule, and $f_{\circ}(\cdot)$ a function that updates the state of the cell at hand as per the states of its neighboring cells. Hence, for a radius $r = 1$ von Neumann’s neighborhood defined over a $d = 2$ -dimensional grid, the set of neighboring cells and state of the cell with coordinates $\mathbf{c} = [i, j]$ are given by:

$$\begin{aligned} f_{\boxplus}([i, j]) &= \{[i, j + 1], [i - 1, j], [i, j - 1], [i + 1, j]\}, \\ S(\mathbf{c}) &= S([i, j]) \\ &= f_{\circ}(S([i, j + 1]), S([i - 1, j]), S([i, j - 1]), S([i + 1, j])), \end{aligned}$$

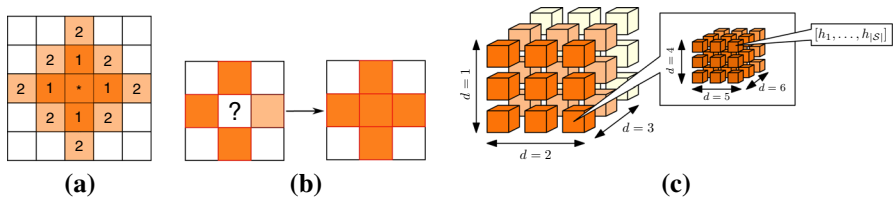


Fig. 1 Neighborhood of CA in data mining: (a) a von Neumann’s neighborhood with radius $r = 1$ and $r = 2$ using the Manhattan distance; (b) the center cell inspects its *von Neumann’s* neighborhood ($r = 1$) and applies the majority voting rule in a one-step update; (c) *CURIE* structure for $d \times \mathcal{G} = 6 \times 3$

i.e., the vector of states $S[i, j]$ of the $[i, j]$ cell within the grid is updated according to the local rule $f_{\odot}(\cdot)$ when applied over its neighbors given by $f_{\boxplus}([i, j])$ (Fig. 1). For a d -dimensional space, a *von Neumann’s* neighborhood contains $2d$ cells.

A cellular automaton should present these three properties: i) *parallelism* or *synchronicity* (all of the updates to the cells compounding the grid are performed at the same time); ii) *locality* (when a cell $[i, j]$ is updated, its state $S[i, j]$ is based on the previous state of the cell and those of its nearest neighbors); and iii) *homogeneity* or *properties-uniformity* (the same update rule $f_{\odot}(\cdot)$ is applied to each cell).

3.2 Cellular automata for data stream mining

A DSM process that may evolve over time can be defined as follows: given a time period $[0, t]$, the historical set of instances can be denoted as $\mathbf{D}_{0,t} = \mathbf{d}_0, \dots, \mathbf{d}_t$, where $\mathbf{d}_i = (\mathbf{X}_i, y_i)$ is an instance, \mathbf{X}_i is the vector of features, and y_i its label. Assuming that $\mathbf{D}_{0,t}$ follows a certain joint probability distribution $P_t(\mathbf{X}, y)$. As it has already been mentioned, data streams usually suffer from *concept drift*, which may change their data distribution, provoking that predictive models trained over them become obsolete. Thus, *concept drift* happens at timestamp $t + 1$ when $P_t(\mathbf{X}, y) \neq P_{t+1}(\mathbf{X}, y)$, i.e. as a change of the joint probability distribution of \mathbf{X} and y at time t .

In addition to the presence of *concept drift*, DSM also imposes by itself its own restrictions, which calls for a redefinition of the previous CA for data mining. Algorithms learning from data streams must operate under a set of restrictions:

- Each instance of the data stream must be processed only once.
- The time to process each instance must be low.
- Only a few data stream instances can be stored (limited memory).
- The trained model must provide a prediction at any time.
- The distribution of the data stream may evolve over time.

Therefore, when adapting a CA for DSM, the above restrictions must be taken into account to yield a CA capable of learning incrementally, and with drift detection and adaptation mechanisms. In order to use CA in DSM, data instances flowing over time must be mapped incrementally to the cells of the grid. Next, we analyze each of the mutually interdependent parts in CA for DSM:

- **Grid:** In a data mining problem with n features, the standard procedure adopted in the literature consists of assigning one grid dimension to each feature. After

that, it is necessary to split each grid dimension by the values of the features, in a way that we obtain the same number of cells per dimension. To achieve that, “bins” must be created for every dimension (Fig. 2) by arranging evenly spaced intervals based on the maximum and minimum values of the features. These “bins” delimit the boundaries of the cells in the grid.

- **States:** We have to define a finite number of discrete states $|\mathcal{S}|$, which will correspond to the number of labels (classes) considered in the data mining problem.
- **Local rule:** In data mining the update rule $f_{\cup}(\cdot)$ can adopt several forms. One of the most accepted variants is a majority vote among neighbors’ states (labels). For example, for $d = 2$:

$$S([i, j]) = \arg \max_{s \in \mathcal{S}} \sum_{[k, l] \in f_{\boxplus}([i, j])} \mathbb{I}(S([k, l]) = s),$$

where the value of $f_{\boxplus}([i, j])$ will be the coordinates of neighboring cells of $[i, j]$, and $\mathbb{I}(\cdot)$ is an auxiliary function taking value 0 when its argument is false and 1 if it is true.

- **Neighborhood:** a neighborhood and its radius must be specified. Even though a diversity of neighborhood relationships has been proposed in the related literature, the “von Neumann” (see Fig. 1) or “Moore” are arguably the most resorted definitions of neighborhood for CA.
- **Initialization:** the grid is seeded with the feature values of the instances that belong to the training dataset. In order to decide the state of each cell, we assign the label corresponding to the majority of training data instances with feature values falling within the range covered by the cell. After that, cells of the grid are organized into regions of similar labels (Fig. 2).
- **Generations:** when the initialization step finishes, some cells may remain unassigned, i.e. not all of them are assigned a state (label). In other words, the training dataset used to prepare the CA for online learning might not be large enough to “fill” all cells in the grid. In such a case, it becomes necessary to “evolve” the grid several times (generations) until all cells are assigned a state. In this evolving process, each cell calculates its new state by applying the update rule over the cells in its neighborhood. All cells apply the same update rule, being updated synchronously and at the same time. Here lies the most distinctive characteristic of CA: the update rule only inspects its neighboring cells, being the processing entirely local (Fig. 1).

4 Proposed approach: CURIE

We delve into the ingredients of *CURIE* to act as drift detector. As shown in Fig. 3, its detection mechanism hinges on the evidence that a recent number of mutations in the neighborhood of a cell that has just mutated, may serve as an symptomatic indicator of the occurrence of a drift.

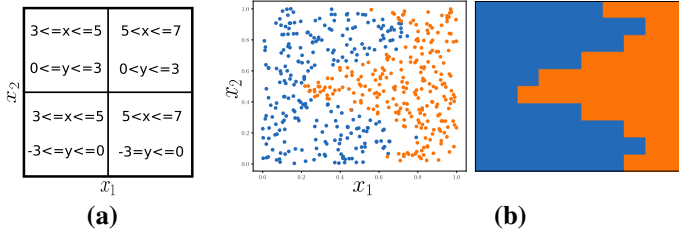


Fig. 2 Data representation in CA: **(a)** a dataset with $d = 2$ dimensions (features), $|S| = \{0, 1\}$, and $\mathcal{G} = 2$ “bins”, where $\mathbf{X}_t = (X_t^1, X_t^2)$ falls between $[3, 7]$ (min/max X_t^1) and $[-3, -3]$ (min/max X_t^2); **(b)** A different dataset whose instances initialize the grid of a bigger cellular automaton with $d = 2$ and $\mathcal{G} = 10$

CURIE builds upon this intuition to efficiently identify drifts in data streams by fulfilling the following methodological steps (Algorithm 1 for the initialization of *CURIE*, and Algorithm 2 for the drift detection and DSM process):

- First, in Algorithm 1 the CA inside *CURIE* is created by setting the value of its parameters (detailed as inputs), and following the characteristics of the given dataset (lines 1–5).
- A reduced number of *preparatory* instances of the data stream $[(\mathbf{X}_t, y_t)]_{t=0}^{P-1}$ is used to initialize the grid of *CURIE*. This grid is seeded with these instances, and then *CURIE* is evolved for several iterations (generations) by applying the local rule until all cells are assigned a state i.e. the labels of the preparatory instances (lines 6–10).
- When the preparatory process is finished, we must ensure that several preparatory data instances have not seeded the same cell, because each cell must reflect only one single state. To this end, we must assign to each cell the most frequent state by inspecting the labels of all those instances that fell within its boundaries. Then, we must ensure that all cells have an assigned state by applying the local rule iteratively all over the grid. Since this last process can again seed a cell with several instances, we have to address this issue to ensure that the cell only reflects one single state (lines 11–13).
- Next in Algorithm 2, *CURIE* starts predicting the data instances coming from the stream in a *test-then-train* fashion (Gama et al. 2014) (lines 2–16). This process consists of first predicting the label of the incoming instance, and next updating the limits of the cells in the grid should any feature value of the processed instance fall beyond the prevailing boundaries of the grid (lines 4–6). Secondly, the label of the incoming instance is used for training, i.e. for updating the state of the corresponding cell (line 7).
- In line 3 *CURIE* stores the incoming instance in a sliding window W of size P , which is assumed, as in the related literature, to be small enough not to compromise the computational efficiency of the overall approach.
- During the *test-then-train* process, *CURIE* checks if a *mutation* of the cell states has occurred (line 9). If the previous state of the cell (before the arrival of the incoming instance) is different from the label of the incoming instance, a mutation has happened. When there is a mutation, we assign the current time step to the cell

in the grid of time steps (line 10). Then, *CURIE* checks the state of the neighboring cells in a radius r_{mut} (of a von Neumann’s neighborhood) in a specific period of time (line 11). If the number of neighboring mutants exceeds a threshold (line 12), *CURIE* considers that a drift has occurred.

- After drift detection, it is time to adapt *CURIE* to the detected change in the stream distribution. To this end, we reset the grid, the vector of states, and the vector of time steps in which a mutation was present (lines 13–15). Finally, the preparatory process is carried out by seeding the grid with the instances stored in the sliding window W (line 16).

Algorithm 1: Steps for the initialization of *CURIE*.

Input : $[(\mathbf{X}_t, y_t)]_{t=0}^{P-1}$; \mathcal{G} ; $f_{\circ}(\cdot)$; $f_{\boxplus}(\mathbf{c})$ for cell with coordinates $\mathbf{c} \in \mathcal{G} = \{1, \dots, G\}^d$.
Output: *CURIE* initialized.

- 1 Let d be equal to the number of features in \mathbf{X}_t
 - 2 Let $|S|$ be the number of classes (alphabet of y_t)
 - 3 Set a vector of state hits per cell: $\mathbf{h}_{\mathbf{c}} = [] \forall \mathbf{c} \in \mathcal{G}$
 - 4 Initialize the limits of the grid: $[(lim_n^{low}, lim_n^{high})]_{n=1}^d$
 - 5 Create the grid as per \mathcal{G} , n and $[(lim_n^{low}, lim_n^{high})]_{n=1}^d$
 - 6 **for** $t = 0$ **to** $P - 1$ **do** // Preparatory process
 - 7 Update limits as per \mathbf{X}_t , e.g., $lim_n^{low} = \min\{lim_n^{low}, x_t^n\}$
 - 8 Update grid “bins” as per \mathcal{G} and $[(lim_n^{low}, lim_n^{high})]_{n=1}^d$
 - 9 Select the cell \mathbf{c} in the grid that encloses \mathbf{X}_t
 - 10 Append y_t to the vector of state hits $\mathbf{h}_{\mathbf{c}'} = [\mathbf{h}_{\mathbf{c}'}, y_t]$
 - 11 Iterate with r and check $|\mathbf{h}_{\mathbf{c}}|$ to ensure one state per cell in \mathcal{G} (local rule $f_{\circ}(\cdot)$ and $f_{\boxplus}(\mathbf{c})$)
 - 12 Guarantee at least $|\mathbf{h}_{\mathbf{c}}| = 1$ in all cells in \mathcal{G}
 - 13 Iterate with r and recheck $|\mathbf{h}_{\mathbf{c}}|$ to ensure one state per cell (local rule $f_{\circ}(\cdot)$ and $f_{\boxplus}(\mathbf{c})$)
-

Finally, after detailing the ingredients of *CURIE* to act as drift detector, we would like to highlight two improvements over (Lobo et al. 2021) that positively impact on the learning of data distribution:

- If the predicted and the true label do not equal each other, the cell state in *CURIE* is always changed to the class of the incoming instance. Otherwise, if the age of the cell state (T_{age}) was considered, this could impact on drift detection resulting in more detection delay.
- In *CURIE* there is always one state assigned to each cell, thus it is not necessary to check the state of the closest cell among those with assigned state to provide a prediction. The cost of assigning one state to all cells of the grid is insubstantial; it is just carried out at the preparatory process and when drift is detected. And so we achieve a more simple method that does not need to check the surroundings (neighborhood) of the cell when no state is assigned.

The source code of *CURIE* is available at <https://github.com/TxusLopez/CURIE>.

Algorithm 2: Steps of *CURIE* for drift detection and DSM

Input : $\{(\mathbf{X}_t, y_t)\}_{t=1}^{\infty}$; \mathcal{G} ; r ; r_{mut} ; $n_muts_allowed$; $mutation_period$; W of size P ; the limits of the grid $[(lim_n^{low}, lim_n^{high})]_{n=1}^d$; a vector of mutations per time step and cell: $\mathbf{h}_m = []$ $\forall \mathbf{m} \in \mathcal{G}$; a vector of state hits per cell: $\mathbf{h}_c = [] \forall \mathbf{c} \in \mathcal{G}$.

Output: Trained *CURIE* producing predictions $\hat{y}_t \forall t \in [1, P, \infty)$

- 1 Let $|S|$ be the number of classes (alphabet of y_t)
- 2 **for** $t = P$ to ∞ **do** // DSM processing
- 3 Update W with the incoming instance (\mathbf{X}_t, y_t)
- 4 Predict \hat{y}_t as $S(\mathbf{c})$, with \mathbf{c} denoting the coordinates of the cell enclosing \mathbf{X}_t
- 5 Update limits as per \mathbf{X}_t , e.g., $lim_n^{low} = \min\{lim_n^{low}, x_t^n\}$
- 6 Update “bins” as per \mathcal{G} and $[(lim_n^{low}, lim_n^{high})]_{n=1}^d$
- 7 Save the current cell state: $cur_st = S(\mathbf{c})$
- 8 Update $S(\mathbf{c}) = y_t$ (i.e. the verified class of test instance)
- 9 **if** $cur_st \neq y_t$ **then** // A mutation occurs in cell
- 10 Append t to the vector of mutations: $\mathbf{h}_{m'} = [\mathbf{h}_{m'}, t]$
- 11 Calculate # mutant neighbors n_muts of the cell, within radius r_{mut} and within time $mutation_period$
- 12 **if** $n_muts \geq n_muts_allowed$ **then** // Detection
- 13 Initialize $\mathbf{h}_m, \mathbf{h}_c$
- 14 Initialize grid limits: $[(lim_n^{low}, lim_n^{high})]_{n=1}^d$
- 15 New grid as per $\mathcal{G}, n, [(lim_n^{low}, lim_n^{high})]_{n=1}^d$
- 16 Preparatory process (lines 6 – 10 of Algorithm 1) with instances in W

5 Experimental setup

In order to assess the performance of *CURIE*, we have designed several experiments with synthetic datasets configured with both abrupt and gradual drift versions.

Since drift detectors usually base their detection mechanisms on the prediction results of a base learner, both detection and classification are often set to work together. As it has been already mentioned, *CURIE* does not use the prediction of the base learner. Instead, it estimates the occurrence of the drift by looking at the changes that occur in the neighborhood of cells deployed over a grid that represents the distribution of data. In our experiments we have accordingly combined three well-known base learners (HT, NB and KNN) with five popular drift detectors including our proposed detector (corr. DDM, EDDM, PH, ADWIN, and *CURIE*). They form 15 different learning-detection schemes following the algorithmic template shown in Algorithm 3. Such base learners and drift detection methods have been selected due to their wide use by the DSM community, and the availability of their implementations in the scikit-multiflow¹ framework. For more information, we refer the reader to Gonçalves Jr et al. (2014) and Barros and Santos (2018). Please note that the inclusion of KNN is not only based on its frequent use, and it has also been considered due to its similarities with CA. While KNN is not strictly local (the neighborhood is not fixed beforehand and the nearest neighbor of an instance may change), CA has a fixed neighborhood. In CA the local interaction between cells affects the evolution of each cell. We would also like to underline that the size of the sliding window of KNN (*max_window_size* parameter in

¹ <https://scikit-multiflow.github.io/>.

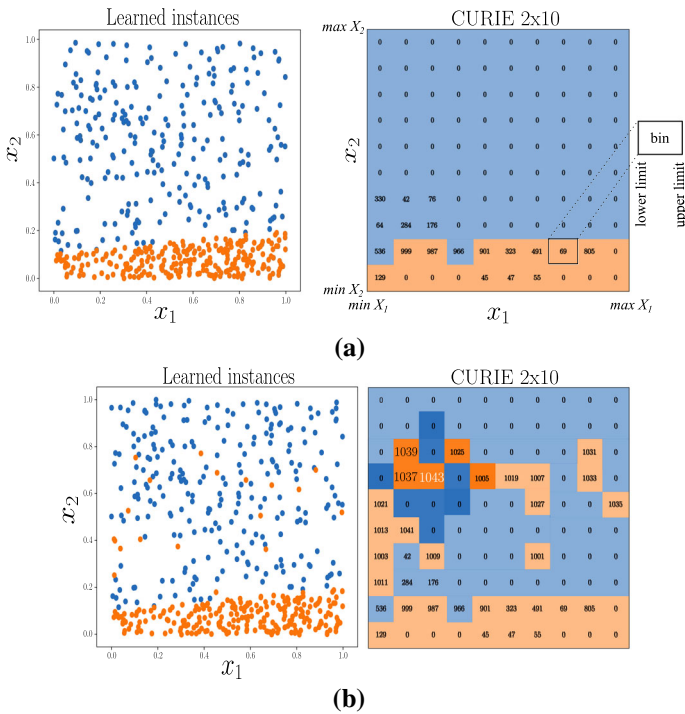


Fig. 3 The interpretable adaptation mechanism of *CURIE* ($d \times G = 2 \times 10$) based on the mutations of its neighborhood: **(a)** before the drift. *CURIE* updates the time instants of each mutant cell, i.e. when the previous state of the cell (before the arriving of the incoming instance) is different from the label of the incoming instance itself; **(b)** drift occurs. *CURIE* checks the neighborhood of each cell, and when at least 2 neighboring cells (defined by *n_muts_allowed* parameter) have mutated in the last 10 time steps (as per the *mutation_period* parameter), *CURIE* considers that a drift has occurred. This is what is declared at $t = 1043$ with the cell [2, 6] and its neighborhood of $r = 2$ (Manhattan distance), where 2 of its neighbors have mutated at time steps 1037 and 1039. The number in each cell represents the last time that has mutated. The neighborhood of the cell [2, 6] is darkened, and this cell (highlighted in white) and its neighbors which have mutated ([1, 6], [1, 7]) have been enlarged

Table 1) is the same than the number of recent instances that CA uses to be initialized and seeded after a drift is detected.

The computer used in the experiments is based on a *x86 64* architecture with 8 processors Intel(R) Core(TM) *i7* at 2.70GHz, and 32GB DDR4 memory running at 2, 133 MHz. The source code for the experiments is publicly available at this GitHub repository: <https://github.com/TxusLopez/CURIE>.

5.1 Datasets

In order to assess the performance of a drift detector by measuring the different detection metrics, we need to know beforehand where a real drift occurs. This is only possible with synthetic datasets. The scikit-multiflow framework, as one of the most commonly accepted libraries in stream learning, allows generating a wide variety of

Table 1 Configuration of detectors and CURIE

Detector	Parameters	Value
DDM	$min_num_instances$	30
	α (<i>warning_level</i>)	2.0
	β (<i>out_control_level</i>)	300
EDDM	$min_num_instances$	30
	α (<i>warning_level</i>)	0.95
	β (<i>out_control_level</i>)	0.9
ADWIN	δ	0.002
PH	$min_instances$	30
	δ	0.005
	<i>threshold</i>	50
	α	0.9999
CURIE	$f_{\boxplus}(\cdot)$	von Neumann
	$f_{\circlearrowleft}(\cdot)$	Majority voting
	r, r_{mut}	2,2
	$ S $	{0, 1}
	$d \times \mathcal{G}$	$n_features \times n_bins$
	<i>mutation_period</i>	10
	<i>num_mutants_neighbors</i>	2

synthetic data in order to simulate a realistic occurrence of drifts. The researcher can configure the number of features, type of features (numerical or categorical), the number of samples, the noise impact, the number of classes, etc. to get the synthetic data closer to a real scenario.

Concretely, we have generated 20 diverse synthetic datasets (10 abrupt and 10 gradual) by using several stream generators (Sine, Random Tree, Mixed, Sea, and Stagger) and functions. They have a different number of features, noise, and a balanced binary class. They exhibit 4 concepts and 3 drifts at time steps 10, 000, 20, 000, and 30, 000 in the case of abrupt datasets, and at time steps 9, 500, 20, 000, and 30, 500 in the case of gradual ones. In the latter case, the width of the drift is 1, 000 time steps. All generated data streams have 40, 000 instances in total.

Finally, as it is explained in Sect. 3.2, it is necessary to create “bins” by splitting each grid dimension by the values of the features. For *Sine* and *RT* datasets we have used 20 “bins” per dimension, while for the rest of datasets we have used 10 “bins”. The values have been found experimentally, just knowing that a small grid is not capable of representing the data distribution (e.g. the grid of Fig. 1). Here, we would like to warn other researches by underlining that CURIE exhibits at this moment a drawback that should be considered. Due to its exponential complexity, we recommend the use of CURIE in datasets with a low number of features. This setback can be tackled by carrying out the search over the grid’s cells by parallelizing this process.

The datasets are available at this Harvard Dataverse repository: <https://dataverse.harvard.edu/dataset.xhtml?persistentId=doi:10.7910/DVN/5OWRGB>. And a more

detailed description can be found at: <https://github.com/TxusLopez/CURIE/blob/master/README.md>

5.2 Methods and parameters

As for *CURIE* we have assigned one grid dimension to each feature of the dataset at hand. We also note that we have used a reduced number of instances to warm up the learning-detection schemes, and also *CURIE* (see Algorithms 1 and 2, and 3). The number of instances for this purpose usually depends on the memory or processing time restrictions. In our experiments we have considered a very low number of them in order to simulate a very restrictive real-time environment (see parameter P in Table 1). In all of them, *CURIE* has been configured with a *von Neumann's* neighborhood rather than opting for its *Moore* alternative. A *von Neumann's* neighborhood is linear in the number of dimensions of the instance space, and therefore scales well for problems of high dimensionality. In addition, a *Moore's* neighborhood includes more neighbors, thus we would have to potentially apply the local rule over more cells. This would make the process computationally heavier and less suited for a DSM setting in the preparatory process and after the drift occurs.

The parameters configuration for the drift detectors under consideration was experimentally fixed, and is detailed in Table 1. The number of preparatory instances (P) and the sliding window (W) of size P are shared between *CURIE* and the base learners. Concretely, their values are $P = 50$ and $w = P$. The values for the base learners parameters have been found through a hyper-parameter tuning process (Grid Search with 10 folds in a *Stratified K-Folds* cross-validator using a balanced accuracy metric) carried out with these preparatory instances (see Table 2). For more information about the meaning of the parameters and their values, we refer the reader to <https://scikit-multiflow.readthedocs.io/en/stable/api/api.html>. And to know more about the Grid Search process we refer to https://scikit-learn.org/stable/modules/generated/sklearn.model_selection.GridSearchCV.html. Finally, Algorithm 3 presents the details of the learning and detection scheme followed by the experiments.

5.3 Performance metrics

Regarding the classification accuracy, we have adopted the so-called *prequential accuracy* ($pACC$) (Dawid and Vovk 1999), which is widely applied in streaming classification scenarios. This metric evaluates the base learner performance by quantifying the average accuracy obtained by the prediction of each test instance before its learning in an online *test-then-train* fashion. This accuracy metric can be defined as:

$$pACC(t) = \begin{cases} pACC_{ex}(t), & \text{if } t = t_{ref}; \text{ otherwise} \\ preACC_{ex}(t-1) + \frac{pACC_{ex}(t) - pACC_{ex}(t-1)}{t - t_{ref} + 1}, & \end{cases}$$

Table 2 Grid Search parameters of the base learners

Base learner	Parameters	Values range
KNN	<i>n_neighbors</i>	3, 5, 7, 10, 15
	<i>leaf_size</i>	3, 5, 7, 10, 15
	<i>algorithm</i>	<i>kd_tree</i>
	<i>max_window_size</i>	<i>W</i>
HT	<i>grace_period</i>	25, 75, 150, 300
	<i>tie_threshold</i>	0.001, 0.25075, 0.5005, 0.75025, 1
	<i>split_confidence</i>	0.000000001, 0.025, 0.05, 0.075, 1.0
	<i>split_criterion</i>	'gini', 'info_gain', 'hellinger'
	<i>leaf_prediction</i>	'mc', 'nb', 'nba'
NB	<i>nominal_attributes</i>	None

Algorithm 3: Learning-detection scheme

Input : $[(\mathbf{X}_t, y_t)]_{t=0}^{t=P-1}; [(\mathbf{X}_t, y_t)]_{t=P}^{\infty}$; W of size P

- 1 . **Output**: Trained base learners producing predictions $\hat{y}_t \forall t \in [P, \infty)$
- 2 Base learner \in [HT, NB, KNN]
- 3 Initialize base learners parameters of Table 1
- 4 Detector \in [DDM, EDDM, ADWIN, PH, CURIE]
- 5 Initialize detectors parameters of Table 1
- 6 **for** $t = 0$ to $P - 1$ **do** // Preparatory process
 - 7 | **if** *detector* = *CURIE* **then**
 - 8 | | Train detector with (\mathbf{X}_t, y_t)
 - 9 | | Train base learner with (\mathbf{X}_t, y_t)
- 10 **for** $t = P$ to ∞ **do** // DSM processing
 - 11 | Update W with the incoming instance (\mathbf{X}_t, y_t)
 - 12 | Predict \hat{y}_t
 - 13 | Train base learner with (\mathbf{X}_t, y_t)
 - 14 | **if** *detector* = *CURIE* **then**
 - 15 | | Train detector with (\mathbf{X}_t, y_t)
 - 16 | **else**
 - 17 | | **if** $\hat{y}_t \neq y_t$ **then**
 - 18 | | | *detector.add_element*(0)
 - 19 | | **else**
 - 20 | | | *detector.add_element*(1)
 - 21 | **if** *detector.detected_change*() **then** // Detection
 - 22 | | Initialize detector
 - 23 | | Preparatory process (6 – 7) with instances in W
- 24 Compare classification and detection performance metrics

where $pACC_{ex}(t) = 0$ if the prediction of the test instance at time t before its learning is wrong, and 1 when it is correct. The reference time t_{ref} fixes the first time step used

in the estimation, and allows isolating the computation of the prequential accuracy before and after a drift has occurred.

To know about the resources used by stream learners, we have adopted the measure RAM-Hours proposed in Bifet et al. (2010), based on rental cost options of cloud computing services. Here, 1 RAM-Hour equals 1 GB of RAM dispensed per hour of processing (GB-Hour). In order to analyze the *concept drift* identifications we have used the Matthews correlation coefficient (MCC) as detection metric. MCC is a correlation coefficient between the current and predicted instances. It returns values in the $[-1, 1]$ range, and it is defined as:

$$MCC = \frac{((TP \cdot TN) - (FP \cdot FN))}{\sqrt{(TP + FP) \cdot (TP + FN) \cdot (TN + FP) \cdot (TN + FN)}}.$$

As the high variance may appear in MCC and μD when $TP = 0$ in several datasets due to a) the very different behavior of detectors in such a diverse experimentation, and b) the period enabled to make detections (2% and 10% of the concept size for abrupt and for gradual datasets respectively) is quite restricted and provokes that $TP = 0$ frequently. Thus, we have included the non-detection rate (ND) as a new indicator to be considered. Therefore, ND shows clearly the detection performance without masking the behavior of those detectors which show a competitive MCC but with a high variance. The lower the rate is, the better this detection performance will be. It is defined as:

$$ND = \frac{\text{The number of experiments with TP=0}}{\text{The total number of experiments (30)}}.$$

We have also measured the distance of the drift detection to the real drift occurrence (μD). Finally, it is worth mentioning that the drifts detected within 2% and 10% [for abrupt and gradual drifts) of the concept size after the real drift positions were computed as TP (as in Barros and Santos (2018)], e.g. $TP > 3$ is possible, otherwise they are FP .

5.4 Statistical tests

We have statistically compared the detectors in all datasets by carrying out the Friedman non-parametric statistical test as described in Demšar (2006). This test is the first step to know whether any of the detectors have a performance statistically different (in prequential accuracy, RAM-Hours, μD , MCC , and ND) from the others. The null hypothesis states that all detectors are statistically equal, and in all cases was rejected. Then it is necessary to use a post-hoc test to discover in what detectors there is a statistical difference (in prequential accuracy, RAM-Hours, μD , and MCC), and we used the Nemenyi post-hoc test (Nemenyi 1963) with 95% confidence to compare all the detectors against all the others. The results are graphically presented showing the critical difference (CD) represented by bars and detectors connected by a bar are not statistically different.

6 Results and discussion

In this section we present the mean results for all datasets in Table 3, and the mean results for abrupt and gradual datasets in Tables 4 and 5 respectively. The original results are given in Online Resource 1 and 2, which correspond to the original results for the order functions $F1$ and $F2$ respectively, as it was introduced in Sect. 5.1. At the beginning of these resources, a complete table compiles the results of the experiments. Each detector is hybridized with the base learners and tested over 20 datasets (10 in Online Resource 1 and 10 in Online Resource 2). The whole set of metrics is presented: prequential accuracy ($pACC$), $RAM - Hours$, true positives (TP), false positives (FP), true negatives (TN), false negatives (FN), distance to the drift (μD), precision, recall, MCC , and non-detection rate (ND). The ranking of $pACC$, $RAM - Hours$, μD , MCC , and ND is provided for every experiment. At the end of the resources, 6 summarizing tables are shown. On the left, the “F global rank” shows the mean rank of detectors, while “F global results” shows the mean and standard deviation of the metrics. On the right, these tables (“F abrupt rank”, “F gradual rank”, “F abrupt results”, and “F gradual results”) detail the mean ranks and metrics for abrupt and gradual datasets. Finally, the Online Resource 3 correspond to Tables 3, 4, and 5. Online Resources 1, 2, and 3 can be found as *O_R_1.xlsx*, *O_R_2.xlsx* and *O_R_3.xlsx* respectively at <https://github.com/TxusLopez/CURIE>.

In Table 3 we observe that CURIE achieves the second best $pACC$ metric with 0.828 ± 0.08 , being ADWIN the best. However, CURIE is the worst in terms of $RAM - Hours$ with $9.55 \cdot 10^{-4} \pm 12 \cdot 10^{-4}$. Here, in favor of CURIE, it is worth mentioning that it is competing with well-established detectors whose code has been optimized and tested by the community in the scikit-multiflow framework. Probably, future versions of CURIE will be more competitive in terms of this metric. Regarding detection metrics, CURIE is the best with 303.45 ± 78.06 for μD , 0.37 ± 0.23 for MCC , and 0.20 for ND .

Deepening the types of drifts, we see in Tables 4 and 5 how CURIE shows the best $pACC$ metric (together with ADWIN) with 0.841 ± 0.09 , and the second best with 0.814 ± 0.07 (with DDM immediately after) for abrupt and gradual datasets respectively. In terms of $RAM - Hours$ metric, CURIE exhibits the worst result for abrupt datasets and the second worst for gradual ones. However, CURIE is the best in μD , MCC , and ND , with 118.90 ± 73.80 , 0.17 ± 0.15 , and 0.40 for abrupt datasets. In case of gradual datasets, CURIE is again the best in μD with 488.00 ± 82.32 , the second best in MCC (very close to ADWIN) with 0.57 ± 0.31 , and the best in ND with 0.00 . Both for abrupt and gradual drifts, CURIE shows a very competitive behavior in $pACC$, and leads the ranking in μD , MCC , and ND ; although in $RAM - Hours$ has a poor performance.

We would like to devote particular attention to the results of the ND metric. As it has been already mentioned in Sect. 5.3, ND indicates how effective the detector is. Here, we see that CURIE achieves true detections ($TP > 0$) in the majority of abrupt datasets ($ND = 0.40$), while makes true detections in all gradual datasets ($ND = 0.00$).

According to the mean ranks for all datasets of Fig. 4, ADWIN, CURIE and DDM are the best detectors in terms of $pACC$, yet no statistical differences between them.

Table 3 Mean results and mean ranks of the detectors in each metric for all considered datasets

		DDM	EDDM	ADWIN	PH	CURIE
pACC	score	0.813 ± 0.08	0.790 ± 0.07	0.835 ± 0.07	0.807 ± 0.08	0.828 ± 0.08
	rank	2.72	4.00	2.18	3.24	2.81
RAM-Hours	score	$5.54 \cdot 10^{-4} \pm 6 \cdot 10^{-4}$	$7.50 \cdot 10^{-4} \pm 18 \cdot 10^{-4}$	$5.86 \cdot 10^{-4} \pm 7 \cdot 10^{-4}$	$4.51 \cdot 10^{-4} \pm 6 \cdot 10^{-4}$	$9.55 \cdot 10^{-4} \pm 12 \cdot 10^{-4}$
	rank	3.31	2.56	3.00	2.32	3.82
μD	score	595.00 ± 0.00	482.80 ± 174.76	392.09 ± 189.50	568.68 ± 92.75	303.45 ± 78.06
	rank	3.93	3.22	2.53	3.56	1.76
MCC	score	0.00 ± 0.00	0.06 ± 0.09	0.26 ± 0.25	0.06 ± 0.16	0.37 ± 0.23
	rank	3.93	3.22	2.53	3.56	1.76
ND	score	1.00	0.62	0.45	0.87	0.20
	rank	3.95	2.95	2.57	3.60	1.94

pACC compiles the prequential accuracy results of those base learners (HT, KNN, and NB) which have been hybridized with each detector (DDM, EDDM, ADWIN, PH, and CURIE). RAM - Hours provides the costs of each mentioned hybrid, while μD , ND, and MCC show the results for the detection metrics

Table 4 Mean results and mean ranks of the detectors in each metric for abrupt datasets

		DDM	EDDM	ADMIN	PH	CURIE
pACC	score	0.815 ± 0.07	0.794 ± 0.07	0.841 ± 0.08	0.808 ± 0.08	0.841 ± 0.09
	rank	2.87	4.08	2.12	3.38	2.55
RAM-Hours	score	$4.48 \cdot 10^{-4} \pm 4 \cdot 10^{-4}$	$2.79 \cdot 10^{-4} \pm 2 \cdot 10^{-4}$	$5.29 \cdot 10^{-4} \pm 5 \cdot 10^{-4}$	$4.99 \cdot 10^{-4} \pm 7 \cdot 10^{-4}$	$9.31 \cdot 10^{-4} \pm 13 \cdot 10^{-4}$
	rank	3.10	2.50	3.27	2.33	3.80
μD	score	190.00 ± 0.00	148.00 ± 50.65	148.33 ± 52.92	179.60 ± 21.94	118.90 ± 73.80
	rank	3.88	2.78	2.70	3.57	2.07
MCC	score	0.00 ± 0.00	0.06 ± 0.07	0.14 ± 0.15	0.10 ± 0.22	0.17 ± 0.15
	rank	3.88	3.08	2.52	3.23	2.28
ND	score	1.00	0.53	0.50	0.80	0.40
	rank	3.91	2.66	2.66	3.38	2.39

Note that μD equals 190 (2% of a concept length of 9500) when there are no $T P$ in abrupt datasets; otherwise, if μD would equal i.e. 0, we would favor this metric

Table 5 Mean results and mean ranks of the detectors in each metric for gradual datasets

		DDM	EDDM	ADWIN	PH	CURIE
pACC	score	0.812 ± 0.08	0.786 ± 0.06	0.828 ± 0.07	0.805 ± 0.08	0.814 ± 0.07
	rank	2.57	3.92	2.25	3.10	3.07
RAM-Hours	score	6.60 · 10 ⁻⁴ ± 9 · 10 ⁻⁴	12.21 · 10 ⁻⁴ ± 34 · 10 ⁻⁴	6.44 · 10 ⁻⁴ ± 9 · 10 ⁻⁴	4.03 · 10 ⁻⁴ ± 4 · 10 ⁻⁴	9.80 · 10 ⁻⁴ ± 11 · 10 ⁻⁴
	rank	3.52	2.62	2.73	2.30	3.83
μD	score	1000.00 ± 0.00	817.59 ± 298.87	635.86 ± 326.08	957.77 ± 163.57	488.00 ± 82.32
	rank	3.98	3.35	2.55	3.88	1.23
MCC	score	0.00 ± 0.00	0.6 ± 0.11	0.38 ± 0.35	0.02 ± 0.09	0.57 ± 0.31
	rank	4.03	3.20	2.57	3.93	1.27
ND	score	1.00	0.70	0.40	0.93	0.00
	rank	3.98	3.23	2.48	3.82	1.48

Note that μD equals 1000 (10% of a concept length of 10,000) when there are no TP in gradual datasets; otherwise, if μD would equal i.e. 0, we would favor this metric

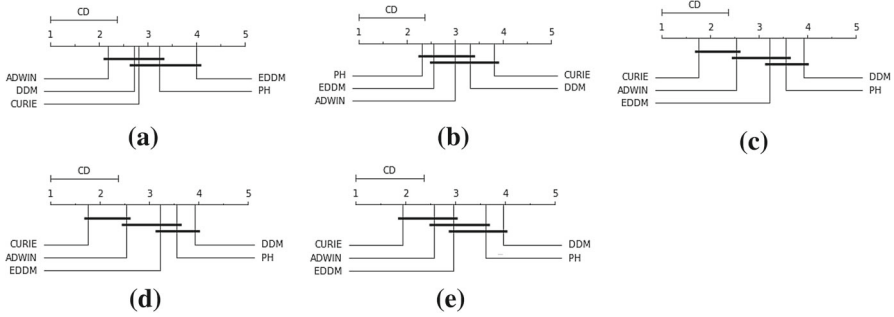


Fig. 4 Comparison of mean ranks for all datasets in (a) pACC, (b) RAM-Hours, (c) μD , (d) MCC, and (e) ND, using the Nemenyi test based on the results of Table 3 with a 95% confidence interval. CD is 1.363887 for 5 detectors and 20 datasets

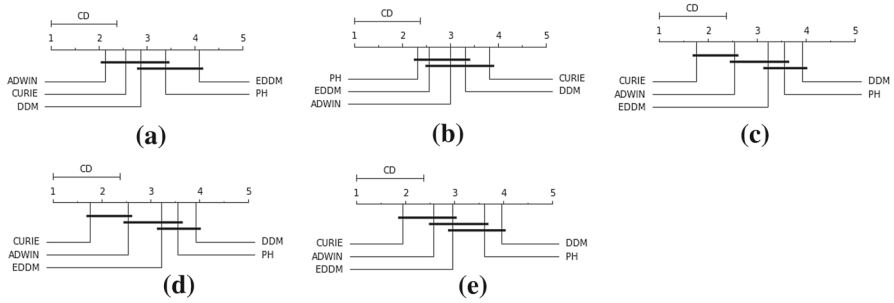


Fig. 5 Comparison of mean ranks for abrupt datasets in (a) pACC, (b) RAM-Hours, (c) μD , (d) MCC, and (e) ND, using the Nemenyi test based on the results of Table 4 with a 95% confidence interval. CD is 1.363887 for 5 detectors and 20 datasets

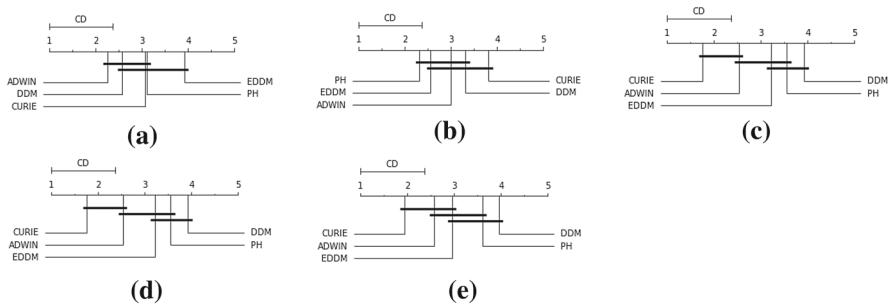


Fig. 6 Comparison of mean ranks for gradual datasets in (a) pACC, (b) RAM-Hours, (c) μD , (d) MCC, and (e) ND, using the Nemenyi test based on the results of Table 5 with a 95% confidence interval. CD is 1.363887 for 5 detectors and 20 datasets

Regarding the *RAM – Hours* metric, CURIE and DDM are the worst detectors, with no statistical differences between them. However, in what refers to μD and MCC, CURIE, ADWIN, and EDDM are the best detectors, yet again no statistical differences between them. We find the same results in the mean ranks for abrupt and gradual datasets depicted in Figs. 5 and 6 respectively.

The conditions of the Nemeneyi test have been very tight (95% confidence for 5 detectors in 20 datasets) and it is difficult to achieve statistical differences. Even so, CURIE has shown to be an interpretable drift detector competitive in terms of predictive performance ($pACC$) and detection metrics (μD , MCC , and ND), without depending on the output (class prediction) of the base learner. Moreover, CURIE provides competitive metrics for abrupt and gradual drifts, being this issue very controversial in the drift detection field, as it was shown in Gonçalves Jr et al. (2014).

Finally, we would like to discuss on the suitability of CA for high-dimensional datasets. When applying CA for streaming scenarios, we have to consider that since the number of cells is proportional to the number of dimensions, as many features in the dataset are, and as many cells per feature, as much computational cost and more processing time will be required. Concretely, in a problem with d dimensions, and with a size of the grid given by G , the worst-case complexity of predicting the class of a given test instance \mathbf{X}_t will be $\mathcal{O}(G^d)$. This is the time required by a single processing thread to explore all cells of the d -dimensional grid and find the cell for the instance \mathbf{X}_t . Because of this exponential complexity, we recommend the use of CA for datasets with a low dimensionality. However, when the scenario imposes such conditions, the search process over the cells of the grid could be parallelized, allowing for faster prediction and cell updating processes.

7 Conclusion and outlook

This work has presented *CURIE*, a competitive and interpretable drift detector based on cellular automata. Until now, cellular automata have shown to be suitable solutions for data mining tasks due to their simplicity to model complex systems, being robust to noise and a low-bias method. Besides, they are computationally complete with parallelism capacity, and they already showed competitive classification performances as data mining methods.

This time, we have focused on their capacity to detect *concept drift*. They have revealed themselves as suitable detectors that achieve competitive detection metrics. They have also allowed base learners to exhibit competitive classification accuracies in a diversity of datasets subject to abrupt and gradual concept drifts. They are suitable candidates to represent data distributions with a few instances, being this ability welcomed in data stream mining tasks where memory and computational resources are often severely constrained. Moreover, *CURIE* can act as an all-in-one approach, in contrast to many other drift detectors which are based on a combination of a base learner method with a detection mechanism.

As future work, we aim to extend the experimental benchmark to more synthetic and real datasets in order to extrapolate the findings and conclusions of this study to different types of drift and more realistic applications. Applying ensemble approaches or even networks of cellular automata are also among our subjects of further study. We also encourage other researchers to improve the performance of CA in high-dimensional datasets. Finally, we would like to underline the relevance of carrying out

experiments on parallelized algorithms; they probably will show the way to overcome the high-dimensionality drawback.

Supplementary Information The online version contains supplementary material available at <https://doi.org/10.1007/s10618-021-00776-2>.

Acknowledgements This work has received funding support from the ECSEL Joint Undertaking (JU) under grant agreement No 783163 (*iDev40* project). The JU receives support from the European Union's Horizon 2020 research and innovation programme, national grants from Austria, Belgium, Germany, Italy, Spain and Romania, as well as the European Structural and Investment Funds. Authors would like to also thank the ELKARTEK and EMAITEK funding programmes of the Basque Government (Spain)

Declarations

Conflict of interest The authors declare that they have no conflict of interest.


References

- Arrieta AB, Díaz-Rodríguez N, Del Ser J, Bennetot A, Tabik S, Barbado A, Salvador G, Sergio GL, Daniel M, Richard B et al (2020) Explainable artificial intelligence (XAI): concepts, taxonomies, opportunities and challenges toward responsible AI. *Inf Fusion* 58:82–115
- Barros RSM, Santos SGTC (2018) A large-scale comparison of concept drift detectors. *Inf Sci* 451:348–370
- Bifet A, Gavalda R, Holmes G, Pfahringer B (2018) Machine Learning for data streams with practical examples in MOA. MIT Press. <https://moa.cms.waikato.ac.nz/book/>
- Bifet A, Holmes G, Pfahringer B, Frank E (2010) Fast perceptron decision tree learning from evolving data streams. In: Pacific-Asia conference on knowledge discovery and data mining. Springer, pp. 299–310
- Carvalho Tiago I, Carneiro Murillo G, Oliveira Gina MB (2019) Improving cellular automata scheduling through dynamics control. *Int J Parallel Emerg Distrib Syst* 34(1):115–141
- Collados-Lara A-J, Pardo-Igúzquiza E, Pulido-Velazquez D (2019) A distributed cellular automata model to simulate potential future impacts of climate change on snow cover area. *Adv Water Resour* 124:106–119
- Dastjerdi AV, Buyya R (2016) Fog computing: helping the Internet of Things realize its potential. *Computer* 49(8):112–116
- Dawid AP, Vovk VG et al (1999) Prequential probability: principles and properties. *Bernoulli* 5(1):125–162
- Del Ser J, Osaba E, Molina D, Yang XS, Salcedo-Sanz S, Camacho D, Das S, Suganthan PN, Coello CAC, Herrera F (2019) Bio-inspired computation: where we stand and what's next. *Swarm Evolut Comput* 48:220–250
- Demšar J (2006) Statistical comparisons of classifiers over multiple data sets. *J Mach Learn Res* 7:1–30
- Fawcett T (2008) Data mining with cellular automata. *ACM SIGKDD Explor News* 10(1):32–39
- Gama J, Žliobaitė I, Bifet A, Pechenizkiy M, Bouchachia A (2014) A survey on concept drift adaptation. *ACM Comput Surv (CSUR)* 46(4):44
- Gilpin W (2019) Cellular automata as convolutional neural networks. *Phys Rev E* 100(3):032402
- Gomes HM, Read J, Bifet A, Barddal JP, Gama J (2019) Machine learning for streaming data: state of the art, challenges, and opportunities. *ACM SIGKDD Explor News* 21(2):6–22
- Gonçalves Jr Paulo M, Santos Silas GT, de Carvalho B, Roberto SM, Vieira Davi CL (2014) A comparative study on concept drift detectors. *Expert Syst Appl* 41(18):8144–8156
- Goumaridis D, Choriantopoulos I, Symeonakis E, Koukoulas S (2019) A random forest-cellular automata modelling approach to explore future land use/cover change in attica (Greece), under different socio-economic realities and scales. *Sci Total Environ* 646:320–335
- Hashemi S, Yang Y, Pourkashani M, Kangavari M (2007) To better handle concept change and noise: a cellular automata approach to data stream classification. In: Australasian joint conference on artificial intelligence. Springer, pp. 669–674

- Hu H, Kantardzic M, Sethi TS (2019) No free lunch theorem for concept drift detection in streaming data classification: a review. In: Wiley interdisciplinary reviews: data mining and knowledge discovery, pp. e1327
- Ilyas M, Mahgoub I (2018) Smart dust: sensor network applications, architecture and design. CRC Press, Boca Raton
- Jie L, Anjin L, Fan D, Feng G, Joao G, Guangquan Z (2018) Learning under concept drift: a review. IEEE Trans Knowl Data Eng 31(12):2346–2363
- Judy JW (2001) Microelectromechanical systems (mems): fabrication, design and applications. Smart Mater Struct 10(6):1115
- Kari J (2005) Theory of cellular automata: a survey. Theor Comput Sci 334(1–3):3–33
- Lobo JL, Del Ser J, Laña I, Bilbao MN, Kasabov N (2018) Drift detection over non-stationary data streams using evolving spiking neural networks. In: International symposium on intelligent and distributed computing. Springer, pp. 82–94
- Lobo JL, Del Ser J, Herrera F (2021) LUNAR: Cellular automata for drifting data streams. Inf Sci 543:467–487
- Losing V, Hammer B, Wersing H (2018) Incremental on-line learning: a review and comparison of state of the art algorithms. Neurocomputing 1275:1261–1274
- Minku Leandro L, Yao X (2011) DDD: a new ensemble approach for dealing with concept drift. IEEE Trans Knowl Data Eng 24(4):619–633
- Nemenyi PB (1963) Distribution-free multiple comparisons. Princeton University, Princeton
- Nichele S, Molund A (2017) Deep learning with cellular automaton-based reservoir computing. Complex Systems
- Pourkashani M, Kangavari MR (2008) A cellular automata approach to detecting concept drift and dealing with noise. In: 2008 IEEE/ACS international conference on computer systems and applications. IEEE, pp. 142–148
- Raghavan R (1993) Cellular automata in pattern recognition. Inf Sci 70(1–2):145–177
- Rudin C (2019) Stop explaining black box machine learning models for high stakes decisions and use interpretable models instead. Nat Mach Intell 1(5):206–215
- Ultsch A (2002) Data mining as an application for artificial life. In: Proceedings of the 5th German workshop on artificial life. Citeseer, pp. 191–197
- Uzun AO, Usta T, Dündar EB, Korkmaz EE (2018) A solution to the classification problem with cellular automata. Pattern Recog Lett 116:114–120
- Von Neumann J, Burks AW et al (1966) Theory of self-reproducing automata. IEEE Trans Neural Netw 5(1):3–14
- Webb GI, Hyde R, Cao H, Nguyen HL, Petitjean F (2016) Characterizing concept drift. Data Min Knowl Disc 30(4):964–994
- Wolfram S (2002) A new kind of science. Wolfram media Champaign, Champaign
- Žliobaitė I, Pechenizkiy M, Gama J (2016) An overview of concept drift applications. In: Big data analysis: new algorithms for a new society. Springer, pp. 91–114

Publisher's Note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

Authors and Affiliations

Jesus L. Lobo¹  · Javier Del Ser² · Eneko Osaba¹ · Albert Bifet³ · Francisco Herrera⁴

✉ Jesus L. Lobo
jesus.lopez@tecnalia.com

Javier Del Ser
javier.delser@ehu.eus ; jdelser@bcamath.org

Eneko Osaba
eneko.osaba@tecnalia.com

Albert Bifet
albert.bifet@telecom-paristech.fr ; abifet@waikato.ac.nz

Francisco Herrera
herrera@decsai.ugr.es

- 1 TECNALIA, Basque Research and Technology Alliance (BRTA), Gipuzkoa, Spain
- 2 University of the Basque Country UPV/EHU & Basque Center for Applied Mathematics (BCAM), Bilbao, Spain
- 3 ParisTech & The University of Waikato, Hamilton, New Zealand
- 4 Andalusian Research Institute in Data Science and Computational Intelligence (DaSCI), University of Granada, Granada, Spain