

Technical Report: Rule-based preprocessing for data stream mining using complex event processing

Aurora Ramírez, Nathalie Moreno,
Antonio Vallecillo

First version: September 2019. Second version: November 2020

Abstract Data preprocessing is known to be essential to produce accurate data from which mining methods are able to extract valuable knowledge. When data constantly arrives from one or more sources, preprocessing techniques need to be adapted to efficiently handle these data streams. To help domain experts to define and execute preprocessing tasks for data streams, this paper proposes the use of active rule-based systems and, more specifically, Complex Event Processing (CEP) languages and engines. The main contribution of our approach is the formulation of preprocessing procedures as event detection rules, expressed in an SQL-like language, that provide domain experts a simple way to manipulate temporal data. This idea is materialised into a publicly available solution that integrates a CEP engine with a library for online data mining. To evaluate our approach, we present three practical scenarios in which CEP rules preprocess data streams with the aim of adding temporal information, transforming features and handling missing values. Experiments show how CEP rules provide an effective language to express preprocessing tasks in a modular and high-level manner, without significant time and memory overheads. The resulting data streams do not only help improving the predictive accuracy of classification algorithms, but also allow reducing the complexity of the decision models and the time needed for learning in some cases.

Contents

1	Introduction	2
2	Background	4
2.1	Data preprocessing	4
2.2	Stream data mining	5
2.3	Complex event processing	6
3	Related work	7
3.1	Combining data mining and CEP	8
3.2	Data stream preprocessing techniques	8
4	CEP for data stream preprocessing	9
4.1	Analysis of CEP functionalities	10
4.2	Implementation	12
5	Experimental evaluation	14
5.1	Methodology	14
5.2	Experiment 1: Generating temporal features	15
5.3	Experiment 2: Integrating categories	19
5.4	Experiment 3: Dealing with missing values	24
6	Discussion	31
6.1	Experimental findings	31
6.2	Comparison with other preprocessing techniques	32
6.3	Strengths and limitations of MOA and CEP preprocessing	36
7	Threats to validity	37
8	Concluding remarks	38

1 Introduction

Data preprocessing is a crucial step in any knowledge discovery process [25]. Cleaning, transforming and integrating data often represent laborious but essential tasks, since patterns are more easily extracted from structured data [61]. Furthermore, the application of an adequate preprocessing method can have a positive impact on the results of the mining process [15,56]. However, data preprocessing is usually a manual error-prone process, which demands some knowledge on statistics and programming skills. Tools like Weka and software packages like those available in R and Python provide support for common preprocessing tasks, but they require a learning curve for non-familiarised domain specialists. Besides, it is likely that they will need to fine-tune parameters or adapt general methods to optimally work on their application domain, making it difficult to reuse procedures.

With the growing interest of organisations in integrating advanced analytics to make informed decisions [51], novel computational techniques to extract knowledge from data streams [22] are highly demanded. The speed at which huge volumes of data are generated pose new challenges to both preprocessing and machine learning (ML) techniques [55]. Despite its relevance, data preprocessing has not been so frequently studied in the context of stream data mining [50]. Most of the available methods try to adapt algorithms applied in traditional data mining or assume that data does not experience any variation over time. Even though specific approaches for dynamic feature selection [5] or discretisation [49] are beginning to appear, they are rarely available in software libraries.

Fast-processing stream engines [1] might represent a candidate platform to implement preprocessing procedures, specially if the input data contain temporal information. Among them, complex event processing (CEP) systems allow users not only to process information flows, but also detect situations of interest and propagate decisions [16]. CEP has experienced a growing interest in the last years due to relevance that Internet of Things (IoT) applications are gaining for, e.g., environmental monitoring [54], healthcare [38], or smart mobility [32].

A distinctive characteristic of CEP systems is that they offer a rich language to express rules and patterns following an SQL-like syntax. For those cases in which transformation and filtering processes depend on the application domain, experts might find it easier to express them in the form of rules. In addition, CEP systems provide specific operators to establish temporal or spatial windows to operate with. This would allow the generation of new features based on a short-term history, which might help to perform more accurate predictions in an incremental learning environment.

Bearing these factors in mind, this paper explores the use of CEP as a novel mechanism to preprocess data streams before the data mining process. In short, the idea underlying our proposal is that CEP characteristics — fast stream processing, rule-based engine and SQL-like syntax — are specially well-suited for domain specialists who might not have deep knowledge of data mining methods and tools. Furthermore, our solution could reduce the burden needed to integrate such mining processes in organisations, thus promoting big data adoption [4]. The development of our proposal is guided by the following research questions (RQ):

RQ1: *How can CEP be used to prepare data streams for online data mining?*

A throughout analysis of the CEP functionalities should be carried out to identify suitable operators and clauses for common preprocessing tasks. The possibility of defining windows to limit the scope of the rules should also be exploited as a way to derive useful temporal information.

RQ2: *Which benefits does CEP bring over other preprocessing methods?* An experimental study is needed to determine whether the resulting data streams have a positive impact on the results of data mining algorithms compared to not applying any preprocessing. Such improvements should be analysed not only in terms of prediction accuracy, but also from other perspectives like model complexity and computational resources (e.g., time and memory).

To the best of our knowledge, this paper represents the first time CEP is applied to define preprocessing rules for data streams. The contributions of this paper are enumerated below:

1. We formulate data stream preprocessing as a rule-based procedure. A theoretical analysis of CEP identifies those functionalities best suited for common preprocessing tasks. As a result, we provide a guide to help domain-specialists choose appropriate operators and functions to define their own stream preprocessing procedures.

2. We present a publicly available implementation of the approach¹ that integrates Esper,² a Java-based CEP engine, with MOA, a library for online data mining [7]. Having both preprocessing and mining steps under the same framework reduces integration efforts and makes information management easier.
3. We conduct three experiments using public datasets from different domains to illustrate the suitability of our proposal in practice. The experiments cover a representative set of preprocessing tasks, after which classification algorithms are applied. The results reveal that the use of the data streams generated by CEP notably improves the prediction performance of the algorithms, reduce the complexity of some decision models and preserve time and memory requirements. These factors are essential to ensure a wider adoption of stream data mining processes in real-world applications.
4. We also discuss the benefits of our approach with respect to other preprocessing alternatives. From this comparison, practitioners gain additional insights of CEP-based preprocessing in terms of native stream characteristics, implementation aspects and adaptability to diverse preprocessing tasks.

The rest of the paper is organised as follows. Section 2 introduces concepts related to data mining and complex event processing. Next, Section 3 presents related work. Our proposal for the use of CEP in data stream preprocessing is explained in Section 4 and experimentally evaluated in Section 5. The strengths and limitations of the approach are discussed in Section 6. Threats to validity are stated in Section 7. Finally, Section 8 concludes and discusses future lines of research.

2 Background

This section explains relevant concepts regarding data preprocessing and stream data mining. Next, the principles underlying CEP are presented.

2.1 Data preprocessing

One of the most important phases within the knowledge discovery from databases (KDD) process consists in preparing the data, a.k.a. instances, from which relevant information has to be extracted. Data preprocessing comprises techniques that adapt either the content or the format of raw data coming from multiple, possibly heterogeneous, sources [25]. The objective is to enhance the accuracy, completeness and consistency of data [29], thus making its management easier in subsequent phases of the KDD process. During the preprocessing step, it is possible to curate the input data by discarding irrelevant or incorrect values, restoring missing values where possible, or detecting outliers. The way

¹ <https://www.github.com/atenearesearchgroup/CEP-Preprocessing>

² Esper v8: <https://www.espertech.com/esper/> (accessed Nov. 27, 2020).

Table 1 Data preprocessing tasks (adapted from [25]).

Data preparation	
<i>Cleaning</i>	Detect and manage missing values and noise.
<i>Transformation</i>	Convert data types and format.
<i>Normalisation</i>	Scale numerical values.
<i>Integration</i>	Join data from different sources.
Data reduction	
<i>Feature selection</i>	Discard irrelevant or redundant features.
<i>Instance selection</i>	Extract data samples, e.g., for validation.
<i>Discretisation</i>	Divide continuous data into intervals.
<i>Value generation</i>	Add new features or generate instances.

in which the attributes, a.k.a. features, of the instances are transformed can largely influence the effectiveness of mining algorithms. When algorithms are provided with quality structured data, they find it easier to detect patterns and produce more reliable predictions in less time [29].

Data preprocessing techniques are usually classified in two main groups according to its purpose: data preparation and data reduction. Table 1 summarises the principal tasks associated to each group [25], which can be applied in combination. Data preparation encompasses methods that detect and fix inconsistencies, e.g., missing values and noise, change data formats or transform values. Sometimes, data coming from multiple sources follow different schema that should be unified too. In contrast, data reduction is focused on altering the number of features or instances to keep only those that satisfy certain criteria. For instance, it is frequent to select subsets of features to reduce the dimensionality of the dataset, whereas some algorithms cannot deal with continuous data, so discretisation is required.

Data preprocessing becomes even more relevant in big data scenarios, where a variety of features and a vast volume of instances are expected [50]. When dealing with data streams, preprocessing techniques should be as lighter and automatic as possible [21]. Additional factors need to be considered due to the dynamic nature of the data. Accurate information regarding its distribution, e.g., ranges of values or dependencies among variables, is not fully available. In the presence of concept drift, the relative importance of the features might change over time [6], making feature selection even harder.

2.2 Stream data mining

A data stream is a potentially unbounded and ordered sequence of instances continuously observed over time [20]. Data stream mining differs from traditional mining processes in a number of aspects [22,46]. From the data perspective, the set of instances is not available beforehand, so they have to be processed one by one or in chunks. Each instance can be accessed a limited number of times (usually only once) and then discarded to optimise the use of memory and storage space. This implies that the decision model is built incrementally, adding new relevant information and forgetting outdated data.

Apart from being able to cope with these issues, stream mining algorithms are expected to provide real-time responsiveness, avoid data queuing and present low memory requirements. Due to the strict conditions in which these algorithms might operate, approximate results are acceptable.

Machine learning methods are popular techniques to carry out the mining step, also when the input data is a stream. Both supervised and unsupervised learning approaches have been developed in the last years [46,31]. Regression and classification are supervised methods that try to predict the value or class of a variable, respectively. In contrast, clustering, whose purpose is to identify groups of related items, and association rule mining, oriented towards extracting patterns describing the data, belong to unsupervised learning.

Focusing on classification, traditional techniques like decision trees need to scan the dataset multiple times. Therefore, the building process has to be adapted to deal with the particularities of data streams. Other algorithms frequently applied in batch learning, such as k-nearest neighbours (kNN) and naive Bayes, do not require substantial modifications as they were originally designed to learn incrementally [22]. Ensemble methods, i.e., those building multiple classifiers, can be designed to learn from data streams too [33], and have become reference techniques due to its good performance and adaptability [13].

The performance of classification algorithms for stream mining can be severely affected when data streams are non-stationary, i.e., the distribution of the attributes or the concept to be predicted changes over time [33]. This phenomenon, known as concept drift, frequently occurs in real-world data streams and can manifest under different patterns, e.g., from gradual shifts to abrupt changes [53]. Dealing with concept drift is an extensive and very active research field, in which different approaches have emerged [13]. Concept drift detectors are algorithms that analyse the stream in order to update the classifier when a change is detected [40,36]. Sliding-window mechanisms create a buffer with samples that is continuously updated. Recent proposals in this area explore the use of multiple windows or dynamically adapt the window size [37, 60]. Finally, classification algorithms can incorporate concept drift detection capabilities, so that they can adapt their learning process [13,28]. Ensemble methods have been also proposed under this approach, updating the role of the base classifiers dynamically [2,14].

2.3 Complex event processing

Complex event processing is a form of information processing aimed at detecting situations of interest from events [41]. In CEP terminology, a *simple* event corresponds to low-level data, e.g., sensor measurements, whereas *complex* events are those derived from simple ones. In order to infer the occurrence of complex events, the expert has to specify rules, which are defined by means of patterns that identify the events to select, and the actions to be accomplished by the rule. Once registered, the CEP engine will analyse the

Table 2 Types of operators in CEP (adapted from [16]).

<i>Selection</i>	Filter events establishing content constraints.
<i>Projection</i>	Extract pieces of information from the events.
<i>Logic</i>	Build conjunctive, disjunctive and negation expressions.
<i>Sequence</i>	Select events using temporal or order conditions.
<i>Flow</i>	Join, divide and sort several event streams.
<i>Creation</i>	Generate new flows and insert events.
<i>Arithmetic</i>	Compute mathematical expressions and statistics.

rules to detect relevant events and accomplish the actions [16]. Traditional CEP systems follow a *publish-subscribe* approach, where subscribers collect rule outcomes to further process them. CEP rules commonly go through three phases:

1. *Selection*. The events triggering the rule are identified within the selected window.
2. *Matching*. The conditions that the event and its attributes should satisfy are checked.
3. *Production*. Complex events are created, possibly computing new attributes with aggregation functions.

Events, rules and patterns are defined using an *event processing language* (EPL), whose syntax is similar to SQL. Listing 1 show the general syntax of a CEP rule in EPL, comprised of three parts: 1) **select**, in which the relevant attributes are indicated (or the whole event, using *); 2) **from**, in which event flows or patterns over them are specified; 3) **where**, in which the conditions to be fulfilled are detailed.

```
select *|<attribute(s)>
from <flow(s)>|<pattern(s)>
where <condition(s)>
```

Listing 1 Syntax of a CEP rule in EPL.

The richness of the EPL language is reflected in its broad range of mathematical, logical and temporal operators, and the flexibility to combine them. Table 2 collects the main types of operators, together with a brief description of their purpose. Another distinctive characteristic of CEP is the definition of *windows* to limit the scope of the rule. A window can be *spatial*, i.e., composed of a number of events, or *temporal*, i.e., time-based defined.

3 Related work

Our work is related to two main lines of research: the combination of data mining and CEP, and the use of data stream preprocessing techniques to improve the results of ML algorithms.

3.1 Combining data mining and CEP

Synergies between CEP and data mining have been explored in the last years, specially in the context of predictive analytics [19, 18]. However, current approaches are mostly focused on applying learning algorithms to enhance CEP capabilities. For instance, ReCEPTor integrates three association rule mining algorithms into a CEP system [47]. The algorithms are defined as new EPL clauses, so that they can be applied by ReCEPTor rules to discover patterns as new events arrive. More recently, Evolving Bayesian Networks have been applied to model changes in the distribution of the event stream processed by CEP [57]. Complex events detected by manually-defined CEP rules have been used as the target to build predictive models [19]. This way, the system is able to anticipate their future appearance.

Given that manually defining CEP rules is a laborious task that largely depends on expert's knowledge, other authors have studied the possibility of automatically inferring such rules. A first proposal, iCEP, breaks down the problem of rule definition into several steps, such as identifying relevant event attributes or determining a proper window size [42]. Different strategies are proposed to solve each task, combining ML techniques like support vector machines (SVM) with information provided by experts. Similarly, autoCEP is a two-phase approach that learns sequences from time series and transforms them into ready-to-deploy CEP rules [45]. Recently, GP4CEP proposes the use of genetic programming to find the conditions that best describe the occurrence of a predefined complex event [11]. In this case, the learned rules are expressed using logical and sequential operators, and can include both temporal and spatial windows. Finally, Adaptive CEP is focused on rule update, for which clustering and Markov probabilistic models are applied [34]. This approach is able to identify similar simple events and then build and dynamically update patterns representing complex events.

3.2 Data stream preprocessing techniques

Several general algorithms to address common data stream preprocessing tasks can be found in the literature. Online methods for normalisation have been recently proposed [9]. For unsupervised learning, the author uses the so-called 'update equations' to estimate the mean and the standard deviation that are employed to normalise the data. For supervised learning, values are scaled using a sigmoid function. Diverse 'update equations' are compared to determine its best scaling factors. The experimental analysis reveals that a simple average estimation performs relatively well. Neither implementation details nor preprocessing times are reported.

Ideas from time series analysis can be used to deal with missing values in numerical features. Recently, the impact of temporal windows on seven imputation methods has been studied in the context of solar irradiance forecasting [17]. Implemented as part of an R package, the compared methods

range from simple statistics, e.g., mean, median and mode, to both linear and nonlinear interpolation. Experiments were conducted injecting missing values at different time rates, a factor that has proven to influence the effectiveness of imputation methods.

Regarding data reduction, dynamic feature selection is probably the most studied topic [6]. Recent approaches rely on different statistics that incrementally determine the relevance of each feature [10,5]. The first method sorts features using the χ^2 metric, counting how many times each category of a feature appears in instances of the same class. The method proposed in the second study to determine the relevance of each feature is considerably more complex. The authors dynamically update a merit function based on entropy that maximises feature relevance while minimising redundancy. Implemented in MOA, the method is evaluated in combination with four classifiers among those available in the library. Interestingly, sometimes the time overhead due to feature selection is compensated with a reduction of the learning time.

Finally, studies for stream discretisation adapt clustering, entropy and frequency-based methods to work in an online setting [24]. A recent proposal is LOFD, a discretisation algorithm that updates interval definition by splitting and merging intervals [49]. The method is included in the MOA extension for data reduction.

4 CEP for data stream preprocessing

This section lays the foundations of a new synergy between CEP and data mining: the use of CEP to build preprocessing rules for data streams. Formally, a preprocessing rule is an expression in the form $C \rightarrow A$, where C is the condition that triggers the rule and A is the action to be performed. For data streams, C represents the arrival of n types of events (E) with different attributes (at): $C = E_1\{at_1, \dots, at_i\} \wedge \dots \wedge E_n\{at_1, \dots, at_j\}$. The action is a set of m operations that transform the events into an instance (I) with new features (f): $A = I\{f_1 = op_1(E_1, \dots, E_n, w(s)), \dots, f_m = op_m(E_1, \dots, E_n, w(t))\}$, where w symbolises an optional window, either spatial (of size s) or temporal (valid for a time interval, t).

These rules can be inspired by the same principles guiding current preprocessing algorithms, but properly adapted to the application domain and expert's knowledge. Furthermore, the use of CEP specific operators and its powerful window handling facilities opens up new possibilities for improving existing data preprocessing procedures and for defining new ones. In response to RQ1, we carry out an analysis of CEP functionalities to identify suitable operators for different kinds of preprocessing tasks. Then, we describe the implementation of the proposed approach, which integrates a CEP engine and incremental ML algorithms.

Table 3 Applicable operators, clauses and functions for each type of preprocessing task.

Data preparation		
Task	Operators	Functions and clauses
<i>Cleaning</i>	- Selection - Logic	<code>any</code> , <code>all</code> , <code>like</code> , <code>some</code> <code>distinctOf</code> , <code>except</code> , <code>exists</code>
<i>Transformation</i>	- Selection - Arithmetic	<code>all</code> , <code>any</code> , <code>like</code> , <code>some</code> <code>cast</code> , <code>Math.round</code> , <code>toDate</code>
<i>Normalisation</i>	- Arithmetic - Logic	<code>avg</code> , <code>maxever</code> , <code>minever</code> , <code>stddev</code> <code>></code> , <code>>=</code> , <code><</code> , <code><=</code>
<i>Integration</i>	- Flow - Sequence - Creation	<code>join</code> , <code>order by</code> , <code>group by</code> <code>first</code> , <code>last</code> , <code>prev</code> , <code>take</code> , <code>→</code> <code>insert into</code> , <code>output</code>
Data reduction		
Task	Operators	Functions and clauses
<i>Feature selection</i>	- Projection	<code>select<attribute></code>
<i>Instance selection</i>	- Selection - Logic - Flow - Sequence	<code>any</code> , <code>all</code> , <code>like</code> , <code>some</code> <code>distinctOf</code> , <code>except</code> <code>group by</code> <code>first</code> , <code>last</code> , <code>prev</code> , <code>take</code> , <code>→</code>
<i>Discretisation</i>	- Selection - Logic	<code>between</code> , <code>in</code> , <code>not in</code> <code>></code> , <code>>=</code> , <code><</code> , <code><=</code>
<i>Value generation</i>	- Arithmetic - Creation - Sequence	<code>avg</code> , <code>min</code> , <code>max</code> , <code>count</code> , <code>sum</code> <code>insert into</code> <code>leastFrequent</code> , <code>mostFrequent</code>

4.1 Analysis of CEP functionalities

Starting from the classification of preprocessing tasks and the types of CEP operators presented in Section 2, we have identified suitable CEP operators for each preprocessing task. Table 3 summarises our analysis, including examples of EPL functions and clauses that implement such operators in Esper. Notice that most of these functions exist in other CEP engines, though their semantic and behaviour might be slightly different.

Regarding data cleaning, CEP rules act as filters that keep only those instances — events, in CEP — satisfying certain conditions. Such constraints should be specified in the **where** clause as logical expressions over event attributes, i.e., the features of the incoming instance. In addition to simple comparison operators (`>`, `<`, `≠`), CEP provides functions like `exists`, which could be applied to detect missing values. More complex expressions can be designed in form of subqueries, using operators like `any` or `some`. Two relevant operators for noise and outlier detection are `distinctOf` and `except`, as they enable setting the values that attributes should present.

Rules for data transformation apply selection operators to identify those event attributes whose values need to be modified. The operators listed in Table 3 are highly flexible, as they work with regular expressions. The arithmetic operator `cast` serves to convert data types, including date formats. Using Esper as CEP engine, it is also possible to include functions from Java packages, e.g., *Math*. If the goal is to normalise values, functions to compute the average

(**avg**) and standard deviation (**stddev**) are available in CEP too. As for maximum values, Esper distinguishes three variants (equivalent for minimum): **max**, to determine the maximum within a window; **fmax**, to include a filter in the expression; and **maxever**, to get the historical maximum.

CEP systems support multiple input streams, which become a relevant feature to tackle data integration during preprocessing. A first approach consists in selecting events from multiple streams using the **from** clause combined with constraints associating related events. This way, it would be possible to unify bank transactions received from different payment systems just by adding a equality comparison on account numbers. The alternative is to employ the **join** clause, which provides a specific syntax to combine streams. CEP includes operators to sort or limit events, combining **output** with **order by** and **limit**, respectively. In addition, events can be redirected to different flows by means of **insert into**. When windows are defined, operators like **first** and **last** are useful to retrieve specific events. Another relevant operator is **followed by** (represented as \rightarrow), since it allows establishing precedence relationships among events. E.g., this operator could be applied to detect purchase habits from which recommendations can be made.

Focusing on data reduction tasks, the **select** clause serves to extract subsets of attributes. This is not actually feature selection as such, since the features of interest have to be known in advance. However, CEP allows nesting **select** clauses, a functionality that could be used to create more complex expressions, e.g., features as result of dynamic queries. Furthermore, it would be possible to define several rules, each one selecting different features, and configure the CEP engine to dynamically activate and deactivate them in response to data distribution changes or temporal conditions.

In contrast, rules for instance selection seem to be easier to materialise, since they allow establishing constraints (**where** clause) on both event content and time of arrival. The combination of selection and logic operators provides the required flexibility to define conditions on attribute values. A simple example here would be a rule that only takes bank transactions exceeding a transfer amount in fraud detection systems. In presence of categorical attributes, the operator **group by** is really effective to split data instances by categories. This would allow specific treatments of data coming from different categories, including their combination. For instance, it would be necessary to create different disease prediction models for men and women.

For discretisation purposes, CEP provides specific operators to express conditions based on numerical intervals. In Esper, intervals are specified by means of the **between** operator as part of the **where** clause. Esper also defines **in** and **not in**, which should be followed by a sequence of values or a query result. Bounds can be established using standard logical operators ($>$, $<$) too. In any case, the ranges could be either set in advance or dynamically adjusted depending on data distribution statistics.

Value generation refers to both adding new features and inserting instances. In the former case, arithmetic operators, e.g., **avg**, **min** and **count**, are included in the **select** clause, requiring the name of the attribute for which they are

computed as parameter. Such functions consider all previous events in the stream or within a particular window, including the current one. As the function is computed for every incoming event, it is actually a new feature that stores temporal information. Adding this type of feature would allow ML algorithms to take short-term history into account, adapting the decision model to data fluctuation. Applications in weather forecasting or those relying on market prices illustrate this point clearly.

Functions inspecting previous events in the flow are relevant for instance generation too. Instead of creating purely artificial instances, a preprocessing method can fill the instance with values obtained by historical statistics. Similarly, frequency-count functions are extremely useful to detect unbalanced data. In such situations, a rule responsible of injecting instances for the minority class should be activated. Previous events might be quite informative to make estimations or identify recurrent values as part of imputation methods to replace missing values.

Most of the aforementioned functions and operators are compatible with the definition of windows to limit their application scope. Furthermore, CEP is characterised by a great parameterisation and nesting capabilities. Finally, many CEP implementations allow users to integrate their own aggregation functions, thus supporting the definition of domain-specific transformations.

4.2 Implementation

Figure 1 shows the high-level structure of the proposed solution, which has been implemented in Java and is available from Github.³ It comprises two main components: the first one (*CEP*) implements data stream preprocessing with CEP, while the second (*DM*) allows invoking data mining algorithms. Each component is explained next in more detail.

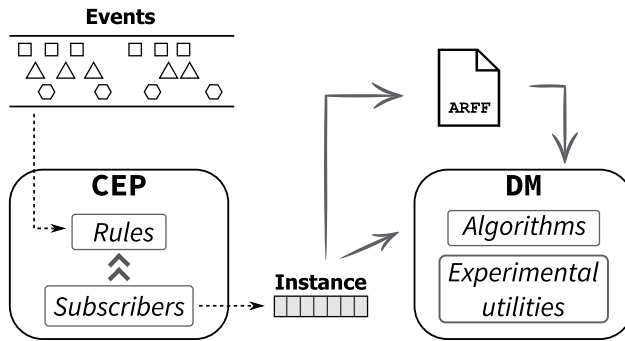


Fig. 1 Core elements of the proposed solution and their relationship.

³ <https://www.github.com/atenearesearchgroup/CEP-Preprocessing>

The *CEP* component receives data in the form of simple events, which might come from multiple streams (represented with different shapes in Figure 1). A number of rules defining the preprocessing actions to be applied should be defined and conveniently registered. Likewise, a subscriber should be associated to each rule, which will be in charge of collecting and processing rule outcomes. From original events and information derived by the rules, the *CEP* component creates the instances to be sent to the incremental learning algorithm. The features of the output instance might represent original event attributes, with or without transformed values, as well as new information (derived attributes). Notice that, if needed, original attributes could be discarded too.

The *DM* component includes the algorithms and utilities required during the data mining phase. Since our proposal is oriented towards real-time data mining, instances are directly sent to the algorithm to train and validate its decision model. Currently available algorithms implement supervised approaches: regression by gradient descent and six classification methods (decision tree, kNN, naive Bayes, a rule-based classifier and two ensemble methods). Finally, it should be noted that the MOA wrapper can be easily extended to invoke other algorithms and utilities from MOA, including drift detectors. For experimental purposes, our implementation also allows loading the generated instances from an ARFF⁴ file, a common format among data mining tools.

Each component has been developed in Java as a *wrapper*. Esper provides the implementation of the CEP engine, whereas the MOA API gives access to data mining algorithms and performance evaluation measures. With the proposed solution, preprocessing a data stream is translated into the following workflow:

1. Define raw data as simple events with attributes in the CEP component (Java object).
2. Define the rule(s) to preprocess the events in EPL language using the operators described in Table 3.
3. Configure the CEP service to register the rules and inject the events in a flow.
4. Process the rule output to create the instance from which the ML algorithm will learn.

Steps 1 and 2 depend on the application domain and therefore should be programmed by the domain specialist. To assist in this process and evaluate the benefits of domain-oriented preprocessing, the next section presents three application scenarios that contribute showing how our approach could be used in practice.

⁴ ARFF format specification: https://waikato.github.io/weka-wiki/formats_and_processing/arff_stable/ (accessed Nov. 27, 2020)

5 Experimental evaluation

This section presents three experiments aimed at illustrating CEP capabilities in representative application domains. The experiments cover different types of preprocessing tasks: feature generation (Section 5.2), data transformation and integration (Section 5.3) and missing value management (Section 5.4). From these experiments, it is possible to analyse the benefits of CEP preprocessing as stated in RQ2. The methodology followed to conduct the experimentation is detailed first.

5.1 Methodology

Input data streams are generated from ARFF datasets publicly available in two well-known repositories: UCI⁵ and OpenML.⁶ Some of these datasets often appear in stream data mining studies [46, 50, 59, 26], so they are representative examples of flows of temporal information that should be processed in real time. The datasets are loaded instance by instance, thus simulating the arrival of simple events to CEP.

A *test-then-train* approach is followed in the learning phase. Each new instance is firstly used to test the current decision model, and then used for training [7]. The decision models are evaluated from the usual perspectives in real-time data mining: prediction performance, time and memory. When applicable, other aspects of the models are discussed.

Four algorithms have been selected for comparative purposes among those available in MOA: *Hoeffding tree* [30], which dynamically induces a decision tree; *k-nearest neighbours* [46], which classifies according to the k most similar instances; *naive Bayes* [46], which predicts the class probability of an instance based on the Bayes theorem; a *rule-based classifier* [23], which derives a set of *if* \rightarrow *then* classification rules; and two ensemble methods with Hoeffding trees as base classifiers, one adopting a bagging strategy (*leveraging bag*) [8] and other including diversity mechanisms and a drift detector (*adaptive random forest*) [27]. These algorithms follow different learning strategies and have been previously used in the literature [48, 49].

To evaluate their performance, we use the window-based evaluator available in MOA with default window size (1000). The following measures are computed: accuracy, i.e., percentage of correctly (positive or negative) classified instances; precision, which represents the percentage of true positive instances among those classified as positive; recall, which returns the percentage of correctly classified instances among all the positive ones; and F_1 , the harmonic mean between precision and recall. The three last measures are computed per class and then reported on average. They all are obtained from the number of true positives (TP), true negatives (TN), false positives (FP) and false negatives (FN) as follows:

⁵ UCI repository: <https://archive.ics.uci.edu/ml> (accessed Nov. 27, 2020)

⁶ OpenML repository: <https://www.openml.org> (accessed Nov. 27, 2020)

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN} \quad (1)$$

$$Precision = \frac{TP}{TP + FP} \quad (2)$$

$$Recall = \frac{TP}{TP + FN} \quad (3)$$

$$F_1 = 2 \cdot \frac{Precision \cdot Recall}{Precision + Recall} \quad (4)$$

Algorithms are executed with and without applying CEP rules, so that the influence of preprocessing decisions, e.g., operators and window size, can be studied. Kruskal-Wallis and Wilcoxon tests ($\alpha = 0.05$) are applied to statistically analyse the results. For pairwise comparisons, p-values are adjusted using the Holm method. As we are interested in analysing the relative performance improvement, all algorithms are executed with default parameters (see the MOA documentation). We consider this would be a common usage scenario for non-experts in data mining. All experiments were run in a Debian 8 computer with 8 cores Intel Core i7-2600 CPU at 3.40 GHz and 16GB RAM. To reduce any possible bias due to machine overload, each experiment was repeated ten times to measure time and memory. Since the two ensemble methods are randomised, a different random seed is configured in each repetition and averaged results are reported. Code, datasets, extended results and statistical validation are available for reproducibility purposes.⁷

5.2 Experiment 1: Generating temporal features

5.2.1 Experiment description

Predicting service demand often require analysing users' consumption habits and price trends. Both aspects are expected to change over time, meaning that any data mining algorithm would probably make better predictions if temporal factors are considered. As a practical example of this situation, a dataset containing the electricity price and demand of two Australian states — New South Wales and Victoria — is studied in this experiment. The dataset,⁸ which is frequently used to evaluate stream data mining algorithms [46], is comprised of 45,312 instances, sampled every 30 minutes. Among the nine features, the dataset contains four numerical features that indicate the price and demand of electricity in each state. The goal is to predict whether the price will rise or drop, so the problem is tackled from a binary classification perspective.

⁷ <https://www.github.com/atenearesearchgroup/CEP-Preprocessing>

⁸ <https://www.openml.org/d/151> (accessed Nov. 27, 2020)

The rationale behind this experiment is twofold. Adding temporal statistics might be relevant to detect the conditions under which price and demand change. In addition, it is necessary to experimentally determine the amount of past information, i.e., window size, that lead to the best performance, since fluctuations of price and demand are not known in advance.

5.2.2 Definition of CEP rules

In order to include the short-term history, we define two CEP rules using aggregation functions over four attributes: electricity price and demand in New South Wales (`nswprice` and `nswdemand`), and equivalent attributes for Victoria state (`vicprice` and `vicdemand`). The first rule (see Listing 2) computes the average value of each of these attributes within a sliding spatial window. Apart from applying the `avg` operator, the rule selects the whole event (*) to keep its original information. Therefore, four derived features are generated. Analogously, the second rule obtains the minimum and maximum values of the attributes. Consequently, eight derived features are created.

```
select  *, avg(nswprice), avg(nswdemand)
         avg(vicprice), avg(vicdemand)
from    Electricity#length(size)
```

Listing 2 Rule that computes the average price and demand of electricity.

We planned two different scenarios to analyse the influence of the derived features. Firstly, the instance is enriched with the temporal features. Secondly, the original price and demand attributes are replaced by the temporal features. Therefore, the CEP engine actually produces four streams: enriched with average values, reduced with average values, enriched with minimum and maximum values, and reduced with minimum and maximum values. In addition, five different window sizes are considered when computing aggregation functions: 10, 50, 100, 500 and 1000.

5.2.3 Results

Table 4 presents the results of the four algorithms in terms of accuracy and F_1 . Symbols (+) and (−) stand for enriched and reduced configurations, respectively. The following acronyms are used for the algorithms: HT (Hoeffding tree), kNN (k-nearest neighbours), NB (naive Bayes), RC (rule-based classifier), LB (leveraging bag) and ARF (adaptive random forest).

Table 4 Accuracy and F_1 results for electricity price change prediction. Figures expressed as percentage, higher values being preferred.

	Window size	Accuracy						F ₁					
		HT	kNN	NB	RC	LB	ARF	HT	kNN	NB	RC	LB	ARF
Original (-)		81.60	82.50	75.30	73.70	92.80	92.20	81.72	82.36	73.30	72.00	92.87	92.31
	10	68.50	76.00	60.80	54.50	77.49	83.26	68.44	75.78	58.33	51.64	77.23	88.06
	50	67.20	73.10	54.70	59.80	75.66	82.20	67.35	72.98	55.11	57.53	75.36	81.85
	100	71.10	72.40	58.00	66.00	75.46	82.27	71.14	72.28	59.09	64.67	75.10	81.90
	500	68.10	71.30	51.60	57.70	75.68	82.34	68.21	71.23	51.65	56.38	75.30	82.00
	1000	67.20	72.50	50.00	57.40	74.19	82.47	66.43	72.32	49.60	55.57	73.86	82.11
	10	80.80	79.90	71.80	71.50	92.87	91.70	80.49	79.88	69.90	70.55	92.90	91.77
	50	98.00	80.20	75.70	84.80	97.98	96.46	97.98	80.05	74.13	85.24	97.98	96.48
Average (+)	100	88.90	80.20	76.90	82.70	92.67	94.66	88.62	79.98	75.43	83.00	92.72	94.73
	500	78.50	78.70	75.60	71.70	91.93	92.50	78.31	78.69	74.03	72.67	91.93	92.61
	1000	84.10	80.10	75.40	63.10	92.29	92.19	83.84	79.91	73.81	60.84	92.29	92.23
	10	74.40	74.40	62.10	62.90	80.13	83.74	74.67	74.14	59.97	60.44	79.62	83.38
Mi/Ma (-)	50	72.20	72.10	58.50	55.90	74.14	81.43	72.52	72.06	59.69	53.87	73.94	81.16
	100	64.40	71.90	59.80	59.30	73.49	80.88	64.67	72.06	61.29	58.16	73.45	80.65
	500	66.40	71.60	50.80	61.30	75.06	82.31	65.75	71.53	50.78	58.95	74.69	82.01
	1000	64.10	71.40	47.90	59.10	75.67	81.61	62.82	71.27	47.33	57.42	75.46	81.25
Mi/Ma (+)	10	80.90	80.50	70.70	82.50	92.60	91.82	80.86	80.43	68.79	81.74	92.62	91.83
	50	91.00	78.40	76.60	65.70	94.19	93.90	91.00	78.40	75.11	63.28	94.27	93.98
	100	88.70	78.40	77.10	78.40	91.56	93.22	88.91	78.46	75.65	77.93	91.71	93.34
	500	85.40	77.10	75.60	72.80	90.84	92.38	85.96	77.09	74.03	72.81	91.01	92.44
1000	82.90	76.50	75.40	77.30	89.89	92.59	82.90	76.46	73.81	76.61	90.09	92.64	

A first interesting finding is that temporal features contribute to improving the performance of most of the algorithms (shaded cells), but only when original features are kept too. This suggests that the current price and demand are relevant for the learning process and should not be omitted. Even so, historical information is highly valuable to complement the decision models, revealing that changes in electricity prices might also be explained by recent consumption habits. Indeed, 30 out of the 60 combinations of algorithm and enriched data stream produce better classification results than using the default data stream.

The suitability of temporal features is specially evident when applying HT, NB and ARF, since the majority of the combinations of aggregation function and window size guarantee better predictions. This fact is a noteworthy achievement considering that ARF already was the second best algorithm among the six tested. Three algorithms with high initial performance (LB, ARF and HT) reach more than 95% accuracy and F_1 under the same preprocessing configuration: average(+) with window size equal to 50. In contrast, no improvement is observed for kNN, a method that classifies according to the k most similar instances. Given that instances closer in time present similar values for the aggregated functions, kNN tend to assign the same class to all of them, not properly capturing the price change due to other features. This issue has a greater impact when the original features are discarded.

In terms of accuracy, the percentage of improvement depends on the applied algorithm: between 1.59% and 20.10% for Hoeffding tree, between 0.13% and 2.39% for naive Bayes, between 4.88% and 15.06% for the rule-based classifier, between 0.08% and 5.58% for leverage bagging and between 0.20% and 4.62% for adaptive random forest. Taking all preprocessed data streams as reference, HT, kNN, LB and ARF outperform NB and RC according to Kruskal-Wallis and two-tailed Wilcoxon tests. As for F_1 measure, most of the algorithms provide a good trade-off between precision and recall. The low values achieved by naive Bayes are due to a marked difference in the recall values for each class. Although this factor is also observed when the original data stream is used, the rate of false positives increases for the configurations with reduced features.

Focusing on preprocessing decisions, the average price and demand seem to be more informative than the minimum and maximum values. Recommended values for the window size are 50 and 100 in both cases, depending on the preferred algorithm. These facts suggest that electricity prices often present fluctuations that are better captured by the average function.

Differences in execution times can be mostly attributed to the selected algorithm, CEP preprocessing only requiring between 0.15 and 0.4 s to process all instances. Almost no difference is observed when the window size is increased, with the exception of the enriched stream with minimum and maximum values. Using kNN increases learning time up to 80 s, since more features necessarily imply more time to measure the distance between instances. For the rest of algorithms, learning time tends to be superior with the preprocessed data streams too, but the achieved performance improvement compensates

this fact, specially when using a fast algorithm like Hoeffding tree. Both this algorithm and naive Bayes are quite stable and always require less than 1 second to conclude. Ensemble methods are considerably more costly in time with either the original stream or those obtained after preprocessing (between 8 and 20 s). The rule-based classifier is the only algorithm able to reduce the learning time a few seconds with respect to the original data stream. Even so, other algorithms are faster and provide better classification performance. As for the memory, running the algorithms with the preprocessed data streams might require nearly double of the memory used with the original data stream, though it never exceeds 11 MB.

It is worth mentioning that some of the decision models derived from the preprocessed data streams are less complex than the original ones. For instance, the number of nodes of the decision tree has been reduced from 57 to 51, while simultaneously increasing accuracy nearly 12% ($min/max(+)$). Furthermore, the highest gain observed in the rule-based classifier (15.6%, $avg(+)$) corresponds to a model with 25 rules instead of the 36 rules originally required. The most discriminatory temporal features, i.e., those most frequently appearing in the decision models, are the average and the maximum electricity prices in New South Wales. The simplicity of the specification of the preprocessing rules using a CEP language, in this case Esper, is also worth noting.

Finally, we visually analyse the positive influence of the added temporal features during the learning phase. Figure 2 shows how accuracy evolves as the data stream is processed. We selected the HT algorithm because it achieves a good trade-off between accuracy and execution time, and reported the highest accuracy (98%). The learning curve when no preprocessing is performed is shown too (dashed line). The addition of the average values of price and demand with a window size equal to 50 provides the best response. The fact that CEP rules operate on a window basis helps capturing changes in the data stream, allowing the algorithm to recover its predictive capability in less time.

5.3 Experiment 2: Integrating categories

5.3.1 Experiment description

Airports constantly schedule and control flights with the aim of providing good service to travellers. Predicting delays and, most importantly, identifying their causes might alleviate inconveniences to users. To address this problem from a ML perspective, the airlines dataset⁹ provides data from 539,383 real flight schedules. This dataset includes four categorical features: the departure and arrival airports, the airline and the day of the week. Airports can take up to 293 distinct values, whereas flights are operated by 18 airlines. The dataset also contains three numerical features: flight identification number (omitted in ML studies), elapsed time and travelling distance. The class attribute indicates whether the flight was delayed (1) or not (0).

⁹ <https://www.openml.org/d/1169> (accessed Nov. 27, 2020)

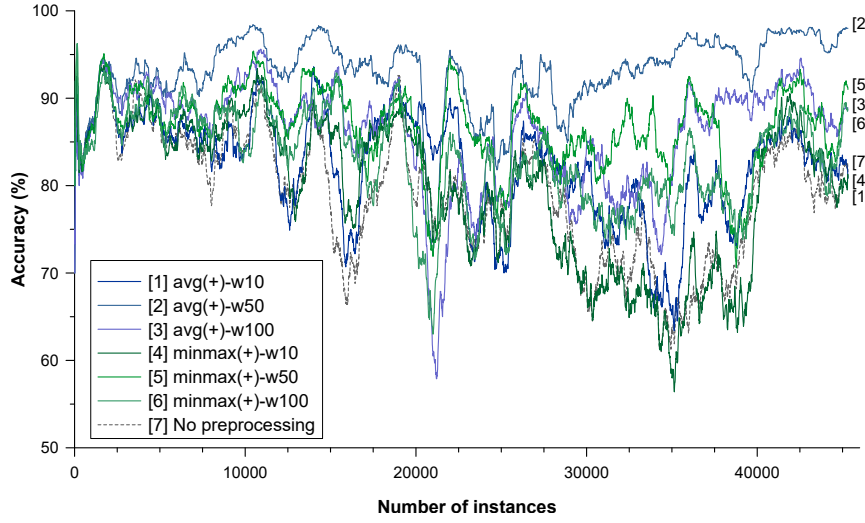


Fig. 2 Learning curve of a Hoeffding tree for electricity data stream enriched with temporal features.

The presence of a high number of categories makes it difficult to generalise decision models. Indeed, a default execution of the Hoeffding tree returns a decision tree with 8,582 nodes and 8,518 leaves, but only four levels of depth. This means that decisions are mostly made on a case-by-case basis, i.e., first splitting by departure airport, then checking the airline and finally deciding depending on the arrival airport. Even though categorical information seems to be highly relevant for learning, this type of decision model is neither manageable nor reusable. Therefore, the purpose of this experiment is to study how categorical information can be indirectly considered in the learning process, in an attempt to not only improve classification performance, but also produce more understandable decision models.

5.3.2 Definition of CEP rules

Two types of rules are defined in this experiment. Firstly, we apply the **group by** clause to aggregate information from the flights according to each categorical feature, i.e., **airportTo**, **airportFrom**, **airline**, and **dayOfWeek**. Listing 3 shows the syntax of the template rule, which selects the original numerical features, i.e., **time** and **length**, as well as the class label (**delay**). Notice that the day of the week is also retrieved, as this feature contains a small number of categories. When the day of week is used in the **group by** clause, the attribute is removed from the **select** clause. We include the **sum** function to count the number of delays within the defined sliding spatial window, but only considering those events that share the same category of the feature appearing in the **group by** clause.

```
select    dayOfWeek, time, length, sum(delay), delay
from      Flight#length(size)
group by feature
```

Listing 3 Rule that computes the number of flight delays by category.

The previous rule groups instances by the category of one feature only, thus allowing us to study the influence of each categorical attribute. Although more than one category could be included in the **group by** clause, the resulting rule would consider those instances sharing all categorical values, which might not be a frequent case if the window size is small. In contrast, the rule detailed in Listing 4 is able to produce an instance with four delay counters, one per categorical feature. It includes subqueries to compute each delay in an independent window, whose name is composed by the prefix ‘w’ and the name of the feature. The **where** clauses match the events in the subquery and the global query. The result of each subquery is renamed using the **as** operator, so that it is easily identified by the subscriber.

```
select  time, length, delay
        (select sum(delay) from Flight#length(size) wAirportF
         where wAirportF.airportFrom=wEvent.airportFrom) as dAirportF,
        (select sum(delay) from Flight#length(size) wAirportT
         where wAirportT.airportTo=wEvent.airportTo) as dAirportTo,
        (select sum(delay) from Flight#length(size) wAirline
         where wAirline.airline=wEvent.airline) as dAirline,
        (select sum(delay) from Flight#length(size) wDayOfWeek
         where wDayOfWeek.dayOfWeek=wEvent.dayOfWeek) as dDayOfWeek,
from      Flight#length(size) wEvent
```

Listing 4 Rule that computes the number of flight delays for all attributes.

5.3.3 Results

Table 5 shows the accuracy and F_1 values obtained for the five data streams produced by the CEP rules. As a baseline, the first row indicates the prediction results for the original data stream.

In view of the high number of shaded cells, relying on delay information of flights sharing certain characteristics leads to a substantial increase of classification performance in all algorithms. Considering either the arrival or the departure airport in a short-term window produces a significant improvement, achieving more than 80% of accuracy with five out of the six algorithms. Slightly better results are obtained with the departure airport (more than 90% of accuracy with all algorithms), suggesting that more delays are explained by the conditions at the origin.

Table 5 Accuracy and F_1 results for flight delay prediction. Figures expressed as percentage, higher values being preferred.

	Window size	Accuracy						F_1					
		HT	kNN	NB	RC	LB	ARF	HT	kNN	NB	RC	LB	ARF
Original	10	64.00	64.60	63.90	45.00	58.86	62.28	62.21	62.29	63.04	51.93	58.57	62.01
	50	84.00	80.50	83.20	62.90	83.58	83.44	85.95	81.02	85.44	57.82	85.72	85.18
	100	65.20	64.40	61.60	56.50	68.23	68.63	69.47	62.46	62.50	52.62	71.40	70.40
	500	63.30	62.80	58.10	63.00	62.37	63.07	64.03	59.91	59.12	63.14	65.01	64.85
	1000	55.40	60.70	49.80	58.50	59.87	56.55	57.41	57.49	51.70	52.59	61.51	59.00
	1000	61.10	62.10	50.10	52.80	57.91	56.44	59.86	58.48	52.63	57.46	59.66	58.62
Airport to	10	96.40	96.40	96.40	90.90	96.40	96.38	96.88	96.88	96.88	90.32	96.88	96.86
	50	87.60	86.70	84.00	65.40	87.59	87.44	89.25	87.81	84.87	60.81	89.24	88.99
	100	80.80	78.30	66.10	70.90	80.75	80.02	81.75	78.86	63.59	70.08	83.09	81.95
	500	60.90	61.60	57.40	59.30	62.59	63.72	63.56	59.25	54.54	62.43	65.98	66.16
	1000	57.40	59.00	54.70	55.00	58.87	58.97	58.07	55.80	52.29	56.59	61.56	61.59
	1000	65.30	61.30	63.60	61.90	65.09	64.22	63.12	58.45	63.91	59.98	63.29	62.61
Day of week	10	60.00	57.20	56.60	54.90	60.10	58.38	58.27	53.38	58.20	57.13	57.62	56.60
	50	58.90	57.40	55.20	58.30	59.59	57.14	56.06	53.75	57.55	59.61	56.88	55.12
	100	58.40	59.00	46.20	53.90	60.29	57.69	57.52	54.60	52.69	57.24	58.23	56.30
	500	54.20	58.40	42.30	52.70	59.55	57.55	55.33	53.86	50.00	57.15	58.45	56.89
	1000	78.90	72.70	62.20	57.70	78.92	78.94	81.68	72.43	56.58	50.00	81.72	81.40
	1000	64.40	63.40	63.10	60.20	64.22	63.68	66.00	61.44	59.88	65.48	66.87	66.70
Airline	10	63.00	62.30	60.60	52.90	62.61	63.10	64.34	59.48	56.87	59.03	64.56	64.54
	500	62.40	61.70	61.20	53.90	60.96	60.32	61.42	58.23	57.17	58.51	61.02	61.05
	1000	62.70	62.30	61.30	56.80	60.62	59.77	60.52	58.50	57.76	58.75	60.64	59.97
	1000	99.20	96.30	92.30	99.40	99.40	99.40	99.31	96.79	93.33	99.48	99.48	99.48
	50	94.80	85.30	72.60	83.20	94.80	94.81	95.49	85.72	73.07	81.31	95.49	95.50
	100	89.50	75.10	65.50	85.90	89.75	89.66	90.74	73.85	65.40	86.14	91.10	90.95
4 counters	500	68.20	64.20	55.80	66.60	68.73	68.72	70.55	61.22	56.30	70.36	70.32	70.15
	1000	62.30	63.50	55.30	61.00	64.56	64.64	64.59	60.26	56.97	63.59	64.25	65.05

In contrast, the day of the week seems to be less informative to predict delays. Accurate information is only gathered if a small number of previous flights is considered, and the achieved improvement — up to 37% — is relatively low compared to airport attributes — up to 102% —. Focusing on the airline, improvements are similar (up to 34%). Furthermore, both accuracy and F_1 results are less dependent on the configured window size. It follows from these results that flights operated by a specific airline are systematically delayed or not, regardless of when and where previous delays occurred.

Even though some improvements are observed when grouping by one feature, combining the derived delay information from the four counters provides the best performance by far. With a window size equal to ten, all algorithms achieve an accuracy greater than 92% (and up to 99% for four algorithms), representing an average percentage of improvement of 66% with respect to the results obtained using the original data stream. F_1 values confirm that high precision and recall are obtained for both positive and negative samples. Ensemble methods, the rule-base classifier and Hoeffding tree stand out as the best methods, nearly reaching perfect classifications. Furthermore, gains are now more generalised across window sizes, kNN and NB being the sole algorithms for which improvement is not possible in two cases. The one-tailed Wilcoxon test confirms that RC, LB and ARF perform better when the four delay counters — with any window size — replace categorical features in the input data stream.

In terms of preprocessing time, CEP rules that group delays by one feature run in less than 2 s, not being affected by the window size. However, the rule producing the four counters requires between 3 and 30 s to conclude, depending on the window size. Additionally, some differences among algorithms are perceived. Hoeffding tree and naive Bayes perform the learning phase with the original data stream in around 5 and 4 s, respectively. Both algorithms become faster after transforming the data stream, requiring less than 3 s to conclude when one counter is applied. Note that the total execution time is less than that of the original stream even if preprocessing time is added. In contrast, the rule-based classifier, which initially required around 10 s, suffers from the presence of numerical attributes. Now it takes several minutes to compute the internal statistics¹⁰ to decide which values increase the predictive performance the most. No great differences are observed for kNN, since categorical values are internally treated as numbers in MOA. Only a slightly increase is perceived for the configuration with four counters. kNN requires around five minutes to execute in all cases, including training with the original data stream. Lastly, ensemble methods (LB and ARF) significantly reduce the time needed for learning if preprocessing is enabled. These algorithms initially required more than 10 minutes to conclude, dropping to less than 3 minutes after using CEP preprocessing.

Regarding memory requirements, it is worth mentioning that all algorithms consume less memory after preprocessing with the exception of RC. The most

¹⁰ Note that this method temporarily stores such statistics in a text file.

significant reductions correspond to the Hoeffding tree method, from nearly 16 to 0.4 MB (on average). The increase of memory of the rule-based classifier, from 3 to 8 MB, is attributed to the growth of the data structures due to a higher number of possible comparisons to build the rules. Despite this, the total memory consumption is still acceptable.

Differences regarding the characteristics of the decision models should be highlighted. The structure of the decision tree — four depth levels, 8,582 nodes and 8,518 leaves for the original data stream — has drastically changed. The most accurate tree after preprocessing presents 15 depth levels, 187 nodes and 94 leaves. The most frequently appearing delay counter is the day of week, followed by the airline, the departure airport and the arrival airport. As for the rule-based classifier, the original data stream led to a set of 182 rules, most of them only referring to the departure airport in their antecedent. When the four delay counters are considered, the number of rules decreased down to 50, and they combine information from all attributes. In fact, flight distance and duration, which did not appear neither in the decision tree nor in the rule set, are now frequently included. Consequently, experts are provided with less complex models, either in the form of a decision tree or a rule set, that include additional decision factors.

Finally, Figure 3 illustrates the evolution of the accuracy (measured in window intervals) to analyse the impact of changes in the stream. Again, we chose the Hoeffding tree as the learning algorithm, whose original performance is depicted with a dashed line. We can see how both the drop rate and the recovering time can be mitigated after preprocessing with CEP rules. In particular, grouping the flight delays by either departure or arrival airport allow reaching high accuracy rates quickly and, more importantly, maintain them above 70% afterwards. The good results obtained by applying the four counters simultaneously are visible from the beginning, and the accuracy remains fairly constant afterwards. Therefore, it is a robust configuration for online learning.

5.4 Experiment 3: Dealing with missing values

5.4.1 Experiment description

The data streams produced by error-prone sensors are expected to contain missing values that should be properly handled. Alternatives vary from removing incomplete data to designing imputation methods to replace them. The purpose of this experiment is to illustrate how these procedures can be encoded as CEP rules, as well as to analyse their influence in the prediction capabilities of classification algorithms. A common practice to study the performance of imputation methods is to inject missing values following different configurations [3].

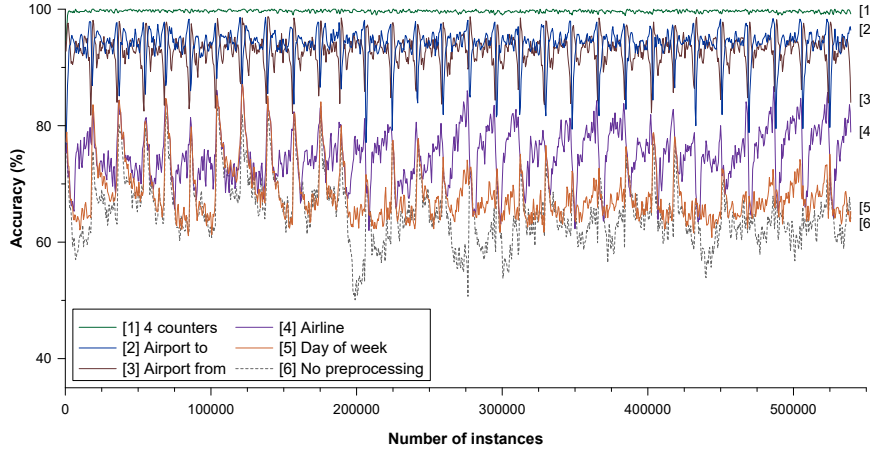


Fig. 3 Learning curve of a Hoeffding tree for airlines data stream with delay counters (window size=10).

For this experiment, we consider a dataset that includes information from light, temperature, humidity and CO2 sensors,¹¹ whose measurements are useful to detect room occupancy [12]. The two samples distributed by the authors are joined to create the input data stream, with 20,560 instances in total. Missing values are inserted following two strategies: random distribution and periodical sequences of missing values. In both cases, we vary the amount of missing values to study its effect on the learning process. For the random strategy, a percentage (5%, 10%, 25%, 50%) of instances is replaced by missing values. For the sequences, both the frequency (100, 500) and the length (5, 10, 15) are configured.

5.4.2 Definition of CEP rules

For this experiment, CEP rules should act differently if the incoming instance contains missing values or not. Complete instances can be detected by means of the rule shown in Listing 5. Assuming that the sensor failure causes a lack of temperature measurement, the temperature attribute (`temp`) would be *null*. Alternatively, other attributes could be included in the **where** clause, or even the absence of the whole event could be detected with the **exists** operator. The instance is simply passed to the learning algorithm.

```
select *
from   Sensor where (temp is not null);
```

Listing 5 Rule that captures complete sensor samples.

¹¹ <https://archive.ics.uci.edu/ml/datasets/Occupancy+Detection+> (accessed Nov. 27, 2020)

If only the previous rule is registered in the CEP engine, instances with missing values will be simply discarded. With the purpose of filling those instances, a second rule should be defined. A possibility is to estimate the missing value considering previous values in the temporal series. We implement an external function to adjust each attribute distribution using the least squares method. When a missing value is detected, the next value in the fitting curve is returned. As for the class label, which specifies whether the room is occupied (1) or not (0), we propose assigning the most frequent value within the window. Listing 6 shows the resulting rule. Notice that the subscriber associated to the rule managing complete instances has to be slightly modified to provide the estimator with samples. In addition, a third rule is periodically triggered to clean out-of-date samples, with the intention of building most accurate fitting curves.

```
select  Function.getEstimation(),
        (select window(occupancy).mostFrequent()
         from Sensor#length(size))
from    Sensor where (temp is null);
```

Listing 6 Rule that generates an instance by estimating missing values.

Two additional imputation methods based on CEP operators are defined for comparative purposes. The first one replicates the last value received in the stream, for which the `lastOf` operator is applied (see Listing 7). As this operator is window-defined, `keepall` is included to specify that the window should contain all previous measurements. The second imputation strategy combines the structure of the two last rules. It computes the average value of each attribute and the class label is set to the most frequent value within a sliding spatial window. For the estimation and the average methods, we evaluate three window sizes: 30, 60 and 90.

```
select
  (select window(temp).lastOf() from Sensor#keepall
   where (temp is not null)),
  (select window(hum).lastOf() from Sensor#keepall
   where (hum is not null)),
  (select window(light).lastOf() from Sensor#keepall
   where (light is not null)),
  (select window(CO2).lastOf() from Sensor#keepall
   where (CO2 is not null)),
  (select window(humRatio).lastOf() from Sensor#keepall
   where (humRatio is not null)),
  (select window(occupancy).lastOf() from Sensor#keepall
   where (occupancy is not null))
from Sensor where (temp is null)
```

Listing 7 Rule that generates an instance by replicating the last available measurement.

5.4.3 Results

Tables 6 and 7 collect the results for data streams with randomly injected missing values (MV) and with missing sequences, respectively. The results for estimation and average methods correspond to the best window size (30). The first row in Table 6 provides the classification performance for the data stream without missing values. In both tables, we also include the accuracy and F_1 obtained when no missing value treatment is applied (referred to as ‘with MV’). It should be noted that, for this experiment, we report absolute values for both performance metrics, since the window-based approach might bias the results if no missing values appear in the window. As it can be observed, all algorithms experience a decrease in their prediction capabilities under the presence of missing values, specially if a high number of randomly distributed instances is missing. These reference configurations allow us to identify the room for improvement of each combination of data stream and algorithm, as well as to analyse whether the imputation method is able to recover the original prediction performance.

According to the one-tailed Wilcoxon test, all algorithms perform better when a curated data stream is received, independently of the imputation method applied. RC slightly decreases in 5 out of the 40 data streams, though the difference in performance is less than 6% lower than the expected value of accuracy or F_1 . Another relevant finding is that classification performance could even be improved with respect to that of the complete data stream. There are two possible reasons explaining this deceptive outcome: the missing instances correspond to false positives or false negatives that are removed, or the imputation method has inverted the class label assigned to such instances. Since the class attribute represents room occupancy, it is expected that positive values appear in sequence, then turning into a sequence of false values, and vice versa. Therefore, the use of the last or the most frequent value within a window could miss the turning point.

Focusing on Table 6, a simple removal of missing values seems to work well. Accuracy returns to values higher than 95% for HT, kNN and NB algorithms. However, estimations based on the average and last values help the rule-based classifier to obtain more robust results, higher than 90%. Ensemble methods (LB and ARF) show similar results, with accuracy values above 98% for all kinds of preprocessing strategies. As it can be expected, imputation methods are more effective as less missing values have to be estimated. Even so, their performance when there is a high percentage of missing values is noteworthy, achieving accuracy percentages only 1% inferior to the ones obtained with the complete data stream. Since samples are taken every minute, consecutive instances in the data stream present similar sensor measurements, so the values returned by window-based estimation methods are quite close to the original ones.

When missing instances appear in sequences, there is no great difference in the final accuracy achieved by the methods (see ‘with MV’ in Table 7). Although all algorithms have time to recover after receiving a missing sequence,

a high rate of such sequences, e.g., 15 missing samples every 100 samples ($\{100,15\}$), degrade their accuracy up to 10% (e.g., from 98.92 to 86.82 for Hoeffding tree). kNN and ensemble methods are less affected in this sense. Again, removing missing values becomes a competitive solution, though accuracy rates closer to the original ones are obtained with imputation methods. In particular, the estimation by the least squares method is highly appropriate to predict next values within the sequence. Methods based on the average and last values are less effective here because they repeat previous measurements. This is reflected in more configurations exceeding the initial accuracy, since more instances share similar values.

Table 6 Accuracy and F_1 results for occupancy prediction (missing values at random). Figures expressed as percentage, higher values being preferred.

	% MV	Accuracy					F ₁						
		HT	kNN	NB	RC	LB	ARF	HT	kNN	NB	RC	LB	ARF
Original		98.92	98.57	96.04	90.84	99.14	99.06	99.03	98.06	96.76	84.75	99.09	98.91
With MV	5	95.23	96.69	95.06	86.75	95.47	95.38	96.57	96.51	94.54	83.91	96.69	96.54
	10	91.69	94.87	94.18	88.86	94.11	91.82	94.30	95.14	92.48	81.45	94.58	94.24
	25	81.94	90.23	91.72	85.40	94.78	91.96	88.00	91.47	86.90	79.99	91.67	93.36
	50	69.04	83.83	88.54	77.21	89.07	93.36	79.71	86.42	78.88	64.37	84.88	93.33
Remove	5	98.94	98.57	96.02	91.45	99.17	99.07	98.97	98.04	96.77	85.93	99.11	98.91
	10	98.48	98.59	96.00	87.96	99.15	99.04	98.73	98.09	96.74	81.83	99.06	98.88
	25	98.16	98.44	95.71	86.71	99.10	99.00	98.43	97.94	96.52	79.89	99.00	98.88
	50	98.82	98.51	96.14	89.34	98.92	98.92	98.86	98.19	96.77	89.63	98.98	98.72
Estimate	5	98.25	98.51	95.72	87.97	99.08	98.98	98.45	97.99	96.43	78.44	99.02	98.79
	10	98.10	98.49	95.57	87.55	98.96	98.93	98.22	97.99	96.34	84.01	98.85	98.75
	25	98.37	98.24	95.00	91.67	98.66	98.65	98.21	97.61	95.79	85.93	98.51	98.37
	50	98.06	98.30	94.40	89.75	98.43	98.49	97.63	97.69	95.10	82.10	98.21	98.14
Average	5	98.91	98.59	96.01	92.72	99.14	99.07	98.95	98.09	96.77	93.95	99.15	98.89
	10	98.89	98.61	96.02	91.33	99.05	99.07	98.92	98.10	96.74	91.52	99.01	98.94
	25	98.74	98.59	95.88	93.72	99.05	99.02	98.55	98.07	96.65	93.27	98.95	98.86
	50	98.63	98.74	95.80	93.57	99.04	99.06	98.62	98.29	96.49	92.98	98.96	98.90
Last	5	98.93	98.59	96.04	90.84	99.13	99.10	99.04	98.10	96.75	84.75	99.10	98.94
	10	98.92	98.54	96.04	90.36	99.12	99.06	99.03	98.03	96.75	84.01	99.05	98.87
	25	98.88	98.52	95.78	90.85	99.11	99.06	98.91	98.04	96.54	89.87	99.03	98.90
	50	98.86	98.46	95.99	90.85	99.10	99.00	98.87	97.97	96.61	84.77	98.97	98.77

Table 7 Accuracy and F_1 results for occupancy prediction (sequences of missing values). Figures expressed as percentage, higher values being preferred.

	Window size	Accuracy						F_1					
		HT	kNN	NB	RC	LB	ARF	HT	kNN	NB	RC	LB	ARF
With MV	{100,5}	95.08	96.61	95.36	88.25	95.43	95.21	96.47	96.54	94.67	82.99	96.61	96.40
	{100,10}	90.82	94.48	94.39	84.23	95.45	91.41	93.73	94.80	92.39	78.71	95.19	93.96
	{100,15}	86.82	92.30	93.42	82.42	92.63	90.71	91.21	93.12	90.20	75.55	92.77	92.74
	{500,5}	98.19	98.22	95.82	89.26	98.40	98.34	98.54	97.73	96.18	81.46	98.61	98.43
	{500,10}	97.46	97.87	95.58	88.89	97.84	97.66	98.06	97.36	95.59	81.07	98.19	97.96
Remove	{500,15}	96.72	97.52	95.36	89.13	98.20	96.97	97.55	97.02	95.02	81.44	97.87	97.45
	{100,5}	98.23	98.63	96.33	92.52	99.12	99.04	98.45	98.14	96.92	88.10	99.05	98.89
	{100,10}	98.39	98.60	96.34	88.04	99.08	99.05	98.62	98.10	96.90	78.11	99.02	98.91
	{100,15}	98.14	98.59	96.34	92.78	99.03	99.05	98.53	98.11	96.97	91.17	98.98	98.84
	{500,5}	98.91	98.55	96.04	88.51	99.12	99.03	99.01	98.02	96.72	84.48	99.06	98.84
Estimate	{500,10}	98.90	98.55	96.03	89.98	99.10	99.04	98.99	98.02	96.67	83.65	99.06	98.86
	{500,15}	98.88	98.54	96.05	91.44	99.08	99.00	98.94	97.98	96.66	86.74	99.04	98.83
	{100,5}	98.19	98.57	96.27	88.21	98.99	99.07	98.33	98.09	96.45	83.74	99.00	98.90
	{100,10}	98.22	98.56	96.14	91.32	98.95	99.04	98.53	98.11	96.02	86.53	98.99	98.89
	{100,15}	98.16	98.61	96.17	93.67	98.89	99.03	98.51	98.26	96.16	92.08	98.98	98.87
Average	{500,5}	98.85	98.53	96.16	90.85	99.11	99.05	98.87	98.02	96.65	84.89	99.05	98.90
	{500,10}	98.81	98.54	96.21	90.58	99.13	99.04	98.94	98.03	96.54	84.71	99.06	98.86
	{500,15}	98.81	98.53	96.24	88.36	99.13	99.05	98.94	98.00	96.58	84.09	99.09	98.86
	{100,5}	98.92	98.64	96.30	91.71	99.11	99.04	98.92	98.17	96.90	91.41	99.02	98.90
	{100,10}	98.20	98.63	96.44	93.39	99.10	99.07	98.39	98.12	97.04	94.39	99.02	98.91
Last	{100,15}	98.16	98.63	96.41	87.82	99.06	99.04	98.41	98.14	97.03	85.96	99.01	98.89
	{500,5}	98.92	98.57	96.00	90.40	99.13	99.06	99.04	98.06	96.66	84.10	99.08	98.89
	{500,10}	98.90	98.53	95.98	88.53	99.14	99.06	99.00	97.98	96.61	83.32	99.10	98.88
	{500,15}	98.86	98.49	95.89	90.90	99.11	99.05	98.98	97.90	96.53	91.19	99.09	98.87
	{100,5}	98.92	98.57	96.37	90.83	99.14	99.05	99.02	98.08	96.97	84.73	99.08	98.88
Last	{100,10}	98.88	98.55	96.43	88.90	99.10	99.04	98.99	98.09	97.00	80.16	99.05	98.85
	{100,15}	98.86	98.56	96.37	88.14	99.08	98.97	98.96	98.10	96.99	79.92	99.06	98.79
	{500,5}	98.92	98.57	96.04	90.84	99.15	99.08	99.03	98.06	96.76	84.75	99.09	98.92
	{500,10}	98.92	98.57	96.04	90.84	99.14	99.05	99.03	98.06	96.76	84.75	99.09	98.89
	{500,15}	98.92	98.57	96.05	90.84	99.15	99.07	99.03	98.06	96.77	84.75	99.09	98.90

6 Discussion

This section provides additional insights from the experimentation. Then, a comparison with other preprocessing techniques is presented. Finally, we also discuss the strengths and limitations that we have observed when using MOA and CEP for data stream preprocessing.

6.1 Experimental findings

To contextualise our experimental results, we have inspected current milestones for the electricity and airlines datasets in OpenML, a platform where researchers can register the outcomes of their ML algorithms. We aim to know whether relatively simple classification algorithms with default parameters, but feeded with preprocessed streams, are competitive compared to state-of-the-art techniques. Looking at the executions collected for the electricity dataset,¹² the best method so far is Ada boost, an ensemble classifier based on decision trees, which achieves 95% of accuracy. Considering that ensemble methods train multiple classifiers, the combination of CEP processing and Hoeffding tree is a competitive alternative for this dataset, since 93% of accuracy is obtained in less than 1 s. Notice that some of the best methods registered in OpenML require more than 5 s to conclude, and might have been trained following an offline approach.

Evaluations of the airlines dataset in OpenML¹³ confirm that it is a challenging dataset, since the best accuracy value currently registered is 67.77%. This result corresponds to the MOA implementation of OzaBag, another ensemble classifier. MOA implementations of kNN, naive Bayes and Hoeffding tree have also been tried with different parameter settings, but all of them learn from the original dataset. Therefore, the choice of the algorithm and its parametrisation seems not to be so relevant compared to producing a high-quality data stream for learning.

However, not all kinds of preprocessing might work for this dataset. For instance, reported accuracy values after discretisation are not higher than 66% for Hoeffding tree and naive Bayes in [49]. These results represent a minor improvement with respect to the values we obtained for the original data stream under our experimental conditions (65.08% and 64.55%, respectively). This suggests that it is rather difficult to improve the accuracy just by modifying numerical features. Indeed, we found that neither flight duration nor distance did appear in the inferred decision tree, which strongly relies on categorical features. The preprocessing proposed here for this data stream allows reducing the presence of categorical information yet keeping its predictive power. As a result, the overall accuracy is increased up to 99% and other features acquire relevance in the decision models. This fact illustrates the need of carefully analysing the problem domain and the behaviour of classifiers to take

¹² <https://www.openml.org/t/219> (accessed Apr. 8, 2020)

¹³ <https://www.openml.org/t/7275> (accessed Apr. 8, 2020)

the most of preprocessing methods. In this respect, counting on languages and tools specific to data stream processing can significantly improve these kinds of analyses, and the effect that different preprocessing alternatives can have on the system.

6.2 Comparison with other preprocessing techniques

To contextualise the results obtained by our rule-based preprocessing approach, this section presents a comparison against other preprocessing techniques. We focus on publicly available techniques coded in the same programming language than our CEP-based solution (Java), which ensures a fair comparison in terms of execution time. Firstly, we have considered the five discretisation algorithms available in MOAReduction: Incremental Discretisation Algorithm (IDA) [58], Incremental Flexible Frequency Discretisation (IFFD) algorithm [39], Local Online Fusion Discretiser (LOFD) [49], Partition Incremental Discretisation (PiD) algorithm [24], and Online ChiMerge (OC) algorithm [35]. These algorithms have been executed with default parameters for the electricity dataset in order to compare their performance against the best configuration of CEP rules reported in Section 5.2.3, i.e., enriching the stream with average price and demand attributes using a window size equal to 50. These discretisation algorithms constitute a representative sample of the state of the art, and have achieved good improvements with this dataset in the past [50]. Secondly, we have applied the MOA filter to replace missing values using the occupancy data streams with injected missing values as input (see Section 5.4.1). This MOA filter has two parameters: the type of replacement for numerical attributes (last, mean, max, min or constant value) and for nominal attributes (last or mode).¹⁴ Combining both parameters, ten different configurations are obtained for comparison. This filter is useful to discuss the results presented in Section 5.4.3, where we proposed four window-based strategies to impute missing values using similar operators. In both comparative studies, Hoeffding tree is chosen as the learning algorithm to be applied after preprocessing. This way, the observed differences in terms of accuracy and execution time can be attributed to the preprocessing step. We compute the same classification performance metrics described in Section 5.1, using a window-based approach for the electricity data stream and the accumulated values for the occupancy data stream.

Table 8 shows the results of the five discretisation algorithms for the data stream used in our first experiment (electricity). The results of the HT algorithm without preprocessing and combined with the best rule-based preprocessing (avg(+) with window size=50) are included for reference too. The last column in Table 8 represents the total execution time (preprocessing and learning) in seconds, reported as the average and standard deviation of ten runs. Focusing on the classification performance, all discretisation algorithms

¹⁴ We omit the option 'Nothing' in both cases, since the resulting data stream is not modified. Default constant value is zero.

Table 8 Classification performance and execution time of discretisation algorithms compared to our best results for electricity dataset.

Preprocessing technique	Accuracy	F_1	Execution time (seconds)
No preprocessing	81.60	81.72	0.44 ± 0.08
IDA	89.00	89.06	56.69 ± 10.22
IFFD	88.20	88.25	54.98 ± 11.57
LOFD	88.80	88.88	57.60 ± 9.58
OC	86.50	86.33	55.87 ± 10.92
PID	87.60	87.53	54.04 ± 12.38
CEP rule (avg(+)-w50)	98.00	97.98	0.89 ± 0.11

are able to increase the accuracy and F_1 values obtained when preprocessing is not applied. The improvements in accuracy range from 6% (OC) to 9% (IDA). However, the proposed rule-based preprocessing approach guarantees a considerable better improvement (20%) than the best discretiser (IDA algorithm). Furthermore, CEP preprocessing has a lower impact on the execution time than the five discretisation algorithms, which can require nearly one minute to conclude. This fact is remarkable considering that MOAReduction uses the same structures than MOA, and the discretisation methods are embedded in the learning process.

Focusing on the occupancy data stream, Figure 4 shows the results of the different imputation methods in terms of accuracy. For each kind of MV injection (random or block), the first column (in grey) represents the accuracy of the Hoeffding tree algorithm without doing any preprocessing. As for the MOA filter, changing the option for nominal attributes (last or mode) does not report any difference in accuracy values, so we only keep the combinations applying the last operator because they are faster. We also rule out the use of the last operator for numerical attributes in the MOA filter, since it never produces an improvement in the learning phase. From the rest of combinations, none of the configurations of the MOA filter is able to recover from the loss of accuracy when only 5% of instances present missing values. In contrast, all CEP preprocessing strategies increase the accuracy values obtained by the learning algorithm. When more missing values appear either randomly or in blocks, the mean operator works better than the rest of options in MOA. Only for two data streams with injected MV (block-100-10 and block-100-15), the MOA filter with mean replacement allows achieving higher accuracy than any of the CEP preprocessing strategies. However, results after CEP preprocessing are more stable, with values always greater than 98% with independence of the amount and distribution of MV. Any CEP operator guarantees better accuracy than the MOA filter with max, min or zero replacement, for which accuracy varies from 91% to 98% depending on the selected operator.

Finally, both preprocessing techniques can be compared in terms of execution time. Table 9 shows the average time (preprocessing and learning) of the ten runs, as well as its standard deviation. Time increased by the MOA filter is quite constant and almost negligible, since learning from the data streams with injected MV took between 0.07 and 0.09 s. CEP preprocessing show more

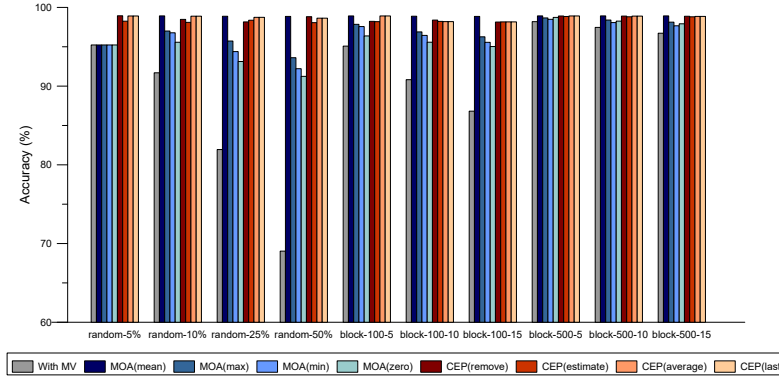


Fig. 4 Percentage of accuracy achieved by Hoeffding tree after applying MOA filter and CEP preprocessing to the occupancy dataset.

variation in this sense, although only the last operator requires more than 1 s to conclude. The difference between MOA filter with mean replacement and CEP preprocessing using the average operator — the two options more similar — can be explained by three factors. Firstly, CEP includes additional operations to convert events into samples, something not required in MOA. Secondly, the CEP engine is event-driven, triggering different rules depending on the incoming event. Doing this distinction is faster in MOA, since all samples are processed by the same class. Thirdly, the CEP operator is configured to use a sliding window, a mechanism not available in MOA at the moment. This last characteristic, together with the possibility of applying a different operators to each attribute, give more flexibility to CEP preprocessing at a minimum cost in execution time.

Table 9 Execution time (in seconds) of Hoeffding tree with different preprocessing techniques to impute MV in the occupancy data stream.

Data stream	MOA(mean)	MOA(max)	MOA(min)	MOA(zero)	CEP(remove)	CEP(estimate)	CEP(average)	CEP(last)
Random (5%)	0.09 ± 2.89E-4	0.10 ± 5.21E-4	0.10 ± 1.28E-3	0.10 ± 1.35E-3	0.17 ± 6.87E-3	0.22 ± 1.49E-2	0.28 ± 9.96E-3	3.25 ± 0.26
Random (10%)	0.09 ± 1.13E-3	0.09 ± 1.19E-3	0.10 ± 1.26E-3	0.10 ± 1.76E-3	0.16 ± 1.80E-3	0.23 ± 1.85E-2	0.29 ± 1.02E-2	5.40 ± 0.54
Random (25%)	0.09 ± 7.99E-4	0.10 ± 4.30E-3	0.11 ± 4.18E-4	0.10 ± 1.29E-3	0.14 ± 1.62E-3	0.24 ± 1.89E-2	0.31 ± 7.94E-3	8.22 ± 0.49
Random (50%)	0.09 ± 3.36E-3	0.11 ± 9.12E-3	0.11 ± 7.75E-3	0.10 ± 5.51E-3	0.12 ± 1.08E-2	0.23 ± 9.44E-3	0.34 ± 1.34E-2	9.80 ± 0.19
Block {100,5}	0.09 ± 4.68E-4	0.10 ± 4.04E-3	0.10 ± 1.42E-3	0.10 ± 2.46E-3	0.17 ± 1.30E-2	0.23 ± 1.65E-2	0.29 ± 1.86E-2	2.39 ± 0.16
Block {100,10}	0.10 ± 1.32E-3	0.10 ± 2.50E-3	0.10 ± 1.76E-3	0.10 ± 8.44E-4	0.16 ± 2.47E-3	0.23 ± 1.45E-2	0.29 ± 1.76E-2	3.39 ± 0.10
Block {100,15}	0.10 ± 1.71E-3	0.10 ± 3.00E-3	0.10 ± 3.74E-4	0.10 ± 7.12E-4	0.16 ± 1.40E-3	0.23 ± 1.95E-2	0.31 ± 2.29E-2	3.87 ± 0.18
Block {500,5}	0.10 ± 3.03E-3	0.10 ± 3.25E-4	0.09 ± 2.20E-3	0.10 ± 1.92E-3	0.17 ± 7.01E-3	0.24 ± 1.95E-2	0.27 ± 7.95E-3	0.94 ± 0.04
Block {500,10}	0.10 ± 1.92E-3	0.10 ± 1.01E-3	0.10 ± 4.48E-4	0.10 ± 9.68E-4	0.18 ± 1.59E-2	0.23 ± 1.57E-2	0.27 ± 9.41E-3	1.21 ± 0.06
Block {500,15}	0.09 ± 3.88E-4	0.10 ± 4.28E-4	0.09 ± 3.54E-4	0.10 ± 1.44E-3	0.17 ± 1.85E-3	0.23 ± 1.73E-2	0.28 ± 1.41E-2	1.56 ± 0.04

6.3 Strengths and limitations of MOA and CEP preprocessing

Table 10 Main characteristics of MOA, MOAReduction and CEP w.r.t. data preprocessing.

	MOA	MOAReduction	CEP
Stream characteristics			
Window	no	partial	yes
Pattern operators	no	no	yes
Multiple inputs	no	no	yes
Needs conversion	no	partial	yes
Implementation details			
Instance	double array	double array	user-defined event
Data types	ARFF types	ARFF types	Java types
Procedure	filter	filter, algorithm	rules
I/O format	ARFF, CSV	ARFF	ARFF, CSV
Preprocessing tasks			
Data preparation	replace MV add noise		replace MV data transformation
Data reduction	attribute selection	discretisation instance selection feature selection	attribute selection feature generation

This section provides an overview of the strengths and limitations of MOA, MOAReduction and CEP, which are summarised in Table 10. As each framework imposes its own processing language, some differences arise regarding how they operate with data streams. In this sense, CEP includes efficient native implementations of concepts that are inherent to stream processing. The support for window definition, both spatial and temporal, as well as the availability of numerous pattern operators, such as **followed by** or **group by** should be highlighted in favour of CEP. These functionalities are not currently supported by MOA, meaning that the user would have to implement them from scratch using Java mechanisms, e.g., queues to simulate windows. The same applies to MOAReduction, though some of its algorithms can be parameterised to specify the frequency of application, a sort of window mechanism.

Another advantage of Esper is the possibility of defining and monitoring multiple input flows, from which rules can be triggered in parallel. Implementing such a functionality in MOA or MOAReduction would require knowledge of concurrency in Java. In contrast, MOA is highly efficient in the management of instances, which are encoded as double arrays. Although MOAReduction follows the same approach, some methods perform instance conversion between the MOA and the Weka implementation of an instance. Conversion is also needed in Esper, since each instance has to be constructed from events implemented as user-defined Java classes. Creating events and converting them to instances can be costly compared to directly handling instances in MOA, as shown in Section 6.2.

Regarding input and output formats, MOA and MOAReduction support ARRF data types, i.e., numerical, nominal, string and date, whereas CEP allows objects and primitive types to be part of the event. In order to use our DM component — built on top of MOA —, output instances have to comply

with the ARFF specification too. However, other tools could be considered for learning purposes after CEP preprocessing, if desired. CSV, another usual format in data analysis, is restricted to clustering in MOA and also presents some limitations in CEP with respect to supported types.

The most distinctive characteristic of our solution is the way preprocessing procedures are implemented. Rules are intuitive mechanisms, but not all preprocessing tasks might be easily expressed in the form of rules. A typical algorithmic flow may be more directly applicable in some situations. For instance, feature selection methods often rely on the dynamic analysis of the classification performance of groups of features, or are embedded in the learning algorithm [50]. At the moment, our approach is designed to apply preprocessing rules independently from the learning process. Contrary to Weka, MOA lacks a preprocessing tab in the GUI, so filters have to be configured as part of the mining process or invoked through code. Likewise, MOAReduction embeds some preprocessing methods into ML algorithms, though independent filters are predefined too.

Focusing on the type of preprocessing tasks currently supported by each framework, some differences also exist. MOAReduction is a specialised extension for data reduction, which covers most of its common tasks (see Table 1). MOA and CEP are more general solutions in this regard, and both can be used to replace missing values and choose attributes. This latter option differs from feature selection, only available in MOAReduction, in that no algorithmic procedure is applied to automatically decide the features to be kept. Regarding value generation and data transformation, MOA only includes a filter to add noise, while discretisation methods in MOAReduction are the only ones that alter numerical features. As for CEP, our experiments have shown how data transformation and feature generation can be addressed by combining CEP operators. Neither MOA nor MOAReduction implement any filter able to enrich the stream with new features in such a flexible way.

7 Threats to validity

This section presents the threats to internal and external validity, and how any possible bias is mitigated.

Internal validity It refers to the aspects of the solution that ensure the causality of the outcomes. A first threat related to how CEP rules are built is the choice of a proper window size. Although our experimental design considers several values for comparison, better values might exist. The selection of ML algorithms constitutes another threat. We opted for publicly available implementations that have previously been considered in the literature [48,49]. Furthermore, they are representative of different categories of algorithms, allowing us to analyse how their internal procedures and decision models are influenced by the preprocessed streams. The use of default parameters might represent an internal threat too. As we are mostly interested in the relative

performance, i.e., before and after applying preprocessing, the default configuration serves our purpose. Moreover, default parameters provide a good baseline according to the results reported in other preprocessing studies and OpenML, as discussed in Section 6.1. Nonetheless, fine-tuned algorithms might achieve better accuracy levels for the original data streams. This may cause a reduction in the improvement attributed to CEP preprocessing, but at the cost of important efforts devoted to parameter tuning.

External validity It mainly concerns the generalisation of the experimental results. Three experiments are conducted to show the applicability of CEP in different preprocessing scenarios. For all of them, several CEP rules are proposed to illustrate the alternatives CEP offers in terms of operators and clauses. The possibilities of applying CEP rules to other preprocessing tasks should be further investigated in the future, though some limitations might exist as pointed out in Section 6.3. For each experiment, a different dataset is used in order to prove that our approach is not limited to a particular application domain. The datasets present different characteristics regarding the number and types of features, as well as number of instances. Some of these datasets often appear in stream data mining studies [46, 50], thus being representative of flows of temporal information that should be processed in real time, e.g., sensor data. Additional datasets are needed to validate our approach in other domains.

8 Concluding remarks

This paper proposes the use of complex event processing as a novel approach to address data preprocessing in the context of stream data mining. Handling raw data as events, the fast-processing CEP engine is able to transform, enrich or curate incoming flows of information making them ready for mining. Such treatments are expressed by means of rules, whose SQL-like syntax helps domain specialists to code them. The analysis of CEP functionalities has revealed a wide range of operators and clauses that fit to different preprocessing purposes, and that can be complemented with the definition of windows and user-coded procedures. Our solution is provided as a Java system that connects Esper, a CEP engine, with MOA, a library for online data mining, thus making it possible to perform preprocessing and learning in one single step.

Three experiments illustrate how rule-based preprocessing can be effectively used for data transformation, feature generation and missing value replacement for diverse application domains. Simple yet powerful preprocessing rules have proven to be able to generate high quality data streams, improving the accuracy of supervised algorithms. In addition, some of the obtained decision models are less complex than those learned from the original data stream, thus making them and the obtained knowledge more manageable and understandable.

These benefits compensate the time and effort required by preprocessing. Our experiments show that time and memory resources do not dramatically increase when new information is added to the learning phase. Also, feeding standard algorithms with CEP-preprocessed streams achieve better results than other types of preprocessing techniques, such as discretisation, in less time.

In the future, we plan to develop more examples to show the suitability of our rule-based approach for other preprocessing tasks, as well as extend its definition to support heterogeneous flows. Our conceptual framework could be implemented in other stream processing engines like Flink or Azure and stream mining solutions like SAMOA [44] and scikit-Multiflow [43]. We hypothesise that the event-based nature of CEP might also be exploited to detect concept drift [52]. More specifically, new types of CEP rules could be defined to analyse the stream data distribution and trigger events depending on the type of concept drift detected. Finally, empirical studies could be designed to analyse whether using a rule-based approach makes data stream preprocessing a lighter task.

Acknowledgements This work was supported by the Spanish Ministry of Economy and Competitiveness under projects TIN2014-52034-R and PGC2018-094905-B-I00.

References

1. Affetti, L., Tommasini, R., Margara, A., Cugola, G., Della Valle, E.: Defining the execution semantics of stream processing engines. *Journal of Big Data* **4**(1), 12 (2017). DOI 10.1186/s40537-017-0072-9
2. Ancy, S., Paulraj, D.: Online Learning Model for Handling Different Concept Drifts Using Diverse Ensemble Classifiers on Evolving Data Streams. *Cybernetics and Systems* **50**(7), 579–608 (2019). DOI 10.1080/01969722.2019.1645996
3. Andiojaya, A., Demirhan, H.: A bagging algorithm for the imputation of missing values in time series. *Expert Systems with Applications* **129**, 10–26 (2019). DOI 10.1016/j.eswa.2019.03.044
4. Baig, M.I., Shuib, L., Yadegaridehkordi, E.: Big data adoption: State of the art and research challenges. *Information Processing & Management* **56**(6), 102,095 (2019). DOI 10.1016/j.ipm.2019.102095
5. Barddal, J.P., Enembreck, F., Gomes, H.M., Bifet, A., Pfahringer, B.: Merit-guided dynamic feature selection filter for data streams. *Expert Systems with Applications* **116**, 227–242 (2019). DOI 10.1016/j.eswa.2018.09.031
6. Barddal, J.P., Gomes, H.M., Enembreck, F., Pfahringer, B.: A survey on feature drift adaptation: Definition, benchmark, challenges and future directions. *Journal of Systems and Software* **127**, 278–294 (2017). DOI 10.1016/j.jss.2016.07.005
7. Bifet, A., Holmes, G., Kirkby, R., Pfahringer, B.: MOA: Massive Online Analysis. *Journal of Machine Learning Research* **11**, 1601–1604 (2010)
8. Bifet, A., Holmes, G., Pfahringer, B.: Leveraging Bagging for Evolving Data Streams. In: J.L. Balcázar, F. Bonchi, A. Gionis, M. Sebag (eds.) *Machine Learning and Knowledge Discovery in Databases*, pp. 135–150. Springer Berlin Heidelberg, Berlin, Heidelberg (2010). DOI 10.1007/978-3-642-15880-3_15
9. Bollegala, D.: Dynamic feature scaling for online learning of binary classifiers. *Knowledge-Based Systems* **129**, 97–105 (2017). DOI 10.1016/j.knosys.2017.05.010

10. Bolon-Canedo, V., Fernández-Francos, D., Peteiro-Barral, D., Alonso-Betanzos, A., Guijarro-Berdiñas, B., Sánchez-Marño, N.: A unified pipeline for online feature selection and classification. *Expert Systems with Applications* **55**, 532–545 (2016). DOI 10.1016/j.eswa.2016.02.035
11. Bruns, R., Dunkel, J., Offel, N.: Learning of complex event processing rules with genetic programming. *Expert Systems with Applications* **129**, 186–199 (2019). DOI 10.1016/j.eswa.2019.04.007
12. Candanedo, L.M., Feldheim, V.: Accurate occupancy detection of an office room from light, temperature, humidity and co2 measurements using statistical learning models. *Energy and Buildings* **112**, 28–39 (2016). DOI 10.1016/j.enbuild.2015.11.071
13. Cano, A., Krawczyk, B.: Evolving rule-based classifiers with genetic programming on gpus for drifting data streams. *Pattern Recognition* **87**, 248–268 (2019). DOI 10.1016/j.patcog.2018.10.024
14. Cano, A., Krawczyk, B.: Kappa Updated Ensemble for drifting data stream mining. *Mach. Learn.* **109**(1), 175–218 (2020). DOI 10.1007/s10994-019-05840-z
15. Crone, S., Lessmann, S., Stahlbock, R.: The impact of preprocessing on data mining: An evaluation of classifier sensitivity in direct marketing. *European Journal of Operational Research* **173**(3), 781–800 (2006). DOI 10.1016/j.ejor.2005.07.023
16. Cugola, G., Margara, A.: Processing Flows of Information: From Data Stream to Complex Event Processing. *ACM Computing Surveys* **44**(3), 15:1–15:62 (2012). DOI 10.1145/2187671.2187677
17. Demirhan, H., Renwick, Z.: Missing value imputation for short to mid-term horizontal solar irradiance data. *Applied Energy* **225**, 998–1012 (2018). DOI 10.1016/j.apenergy.2018.05.054
18. Flouris, I., Giatrakis, N., Deligiannakis, A., Garofalakis, M., Kamp, M., Mock, M.: Issues in complex event processing: Status and prospects in the Big Data era. *Journal of Systems and Software* **127**, 217–236 (2017). DOI 10.1016/j.jss.2016.06.011
19. Fülöp, L.J., Beszédes, A., Tóth, G., Demeter, H., Vidács, L., Farkas, L.: Predictive Complex Event Processing: A Conceptual Framework for Combining Complex Event Processing and Predictive Analytics. In: *Proceedings 5th Balkan Conference in Informatics (BCI)*, pp. 26–31 (2012). DOI 10.1145/2371316.2371323
20. Gaber, M.M.: Advances in data stream mining. *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery* **2**(1), 79–85 (2012). DOI 10.1002/widm.52
21. Gaber, M.M., Zaslavsky, A., Krishnaswamy, S.: Mining Data Streams: A Review. *ACM SIGMOD Record* **34**(2), 18–26 (2005). DOI 10.1145/1083784.1083789
22. Gama, J.: *Knowledge Discovery from Data Streams*, 1st edn. Chapman & Hall/CRC (2010)
23. Gama, J., Kosina, P.: Learning Decision Rules from Data Streams. In: *Proceedings of the 22th International Joint Conference on Artificial Intelligence (IJCAI) - Volume II*, pp. 1255–1260 (2011). DOI 10.5591/978-1-57735-516-8/IJCAI11-213
24. Gama, J., Pinto, C.: Discretization from Data Streams: Applications to Histograms and Data Mining. In: *Proceedings of the 2006 ACM Symposium on Applied Computing (SAC)*, pp. 662–667. ACM (2006). DOI 10.1145/1141277.1141429
25. García, S., Luengo, J., Herrera, F.: *Data Preprocessing in Data Mining*. Springer (2014)
26. Ghomeshi, H., Gaber, M.M., Kovalchuk, Y.: EACD: evolutionary adaptation to concept drifts in data streams. *Data Mining and Knowledge Discovery* **33**(3), 663–694 (2019). DOI 10.1007/s10618-019-00614-6
27. Gomes, H.M., Bifet, A., Read, J., Barddal, J.P., Enembreck, F., Pfahringer, B., Holmes, G., Abdesslem, T.: Adaptive random forests for evolving data stream classification. *Machine Learning* **106**, 1469–1495 (2017). DOI 10.1007/s10994-017-5642-8
28. Hammami, Z., Sayed-Mouchaweh, M., Mouelhi, W., Said, L.B.: Neural networks for online learning of non-stationary data streams: a review and application for smart grids flexibility improvement. *Artif Intell Rev* **53**, 6111–6154 (2020). DOI 10.1007/s10462-020-09844-3
29. Han, J., Kamber, M., Pei, J.: *Data Preprocessing*, 3rd edition edn., chap. 3, pp. 83–124. Morgan Kaufmann (2012). DOI 10.1016/B978-0-12-381479-1.00003-4
30. Hulten, G., Spencer, L., Domingos, P.: Mining Time-changing Data Streams. In: *Proceedings of the 7th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD)*, pp. 97–106 (2001). DOI 10.1145/502512.502529

31. Jin, R., Agrawal, G.: Frequent Pattern Mining in Data Streams, chap. 4, pp. 61–84. Springer US (2007). DOI 10.1007/978-0-387-47534-9_4
32. Kousiouris, G., Akbar, A., Sancho, J., Ta-shma, P., Psychas, A., Kyriazis, D., Varvarigou, T.: An integrated information lifecycle management framework for exploiting social network data to identify dynamic large crowd concentration events in smart cities applications. *Future Generation Computer Systems* **78**, 516–530 (2018). DOI 10.1016/j.future.2017.07.026
33. Krawczyk, B., Minku, L.L., Gama, J., Stefanowski, J., Wozniak, M.: Ensemble learning for data stream analysis: A survey. *Inf. Fusion* **37**, 132–156 (2017). DOI 10.1016/j.inffus.2017.02.004
34. Lee, O.J., Jung, J.E.: Sequence Clustering-based Automated Rule Generation for Adaptive Complex Event Processing. *Future Generation Computer Systems* **66**, 100–109 (2017). DOI 10.1016/j.future.2016.02.011
35. Lehtinen, P., Saarela, M., Elomaa, T.: Online ChiMerge Algorithm, pp. 199–216. Springer Berlin Heidelberg (2012). DOI 10.1007/978-3-642-23241-1_10
36. Liu, A., Song, Y., Zhang, G., Lu, J.: Regional concept drift detection and density synchronized drift adaptation. In: C. Sierra (ed.) *Proc. 26th International Joint Conference on Artificial Intelligence (IJCAI)*, pp. 2280–2286. ijcai.org (2017). DOI 10.24963/ijcai.2017/317
37. Liu, A., Zhang, G., Lu, J.: Fuzzy time windowing for gradual concept drift adaptation. In: *Proc. 2017 IEEE International Conference on Fuzzy Systems (FUZZ-IEEE)*, pp. 1–6. IEEE (2017). DOI 10.1109/FUZZ-IEEE.2017.8015596
38. Loreti, D., Chesani, F., Mello, P., Roffia, L., Antoniazzi, F., Cinotti, T.S., Paolini, G., Masotti, D., Costanzo, A.: Complex reactive event processing for assisted living: The Habitat project case study. *Expert Systems with Applications* **126**, 200–217 (2019). DOI 10.1016/j.eswa.2019.02.025
39. Lu, J., Yang, Y., Webb, G.I.: Incremental discretization for naïve-bayes classifier. In: *Advanced Data Mining and Applications*, pp. 223–238. Springer Berlin Heidelberg (2006)
40. Lu, N., Lu, J., Zhang, G., de Mántaras, R.L.: A concept drift-tolerant case-base editing technique. *Artif. Intell.* **230**, 108–133 (2016). DOI 10.1016/j.artint.2015.09.009. URL <https://doi.org/10.1016/j.artint.2015.09.009>
41. Luckham, D.: *The Power of Events: An Introduction to Complex Event Processing in Distributed Enterprise Systems*. Addison-Wesley (2001)
42. Margara, A., Cugola, G., Tamburrelli, G.: Learning from the Past: Automated Rule Generation for Complex Event Processing. In: *Proceedings of the 8th ACM International Conference on Distributed Event-Based Systems (DEBS)*, pp. 47–58 (2014). DOI 10.1145/2611286.2611289
43. Montiel, J., Read, J., Bifet, A., Abdesslem, T.: Scikit-Multiflow: A Multi-output Streaming Framework. *Journal of Machine Learning Research* **19**(72), 1–5 (2018). URL <http://jmlr.org/papers/v19/18-251.html>
44. Morales, G.D.F., Bifet, A.: SAMOA: Scalable Advanced Massive Online Analysis. *Journal of Machine Learning Research* **16**, 149–153 (2015). URL <http://jmlr.org/papers/v16/morales15a.html>
45. Mousheimish, R., Taher, Y., Zeitouni, K.: Automatic Learning of Predictive CEP Rules: Bridging the Gap Between Data Mining and Complex Event Processing. In: *Proceedings of the 11th ACM International Conference on Distributed and Event-based Systems (DEBS)*, pp. 158–169 (2017). DOI 10.1145/3093742.3093917
46. Nguyen, H.L., Woon, Y.K., Ng, W.K.: A survey on data stream clustering and classification. *Knowledge and Information Systems* **45**(3), 535–569 (2015). DOI 10.1007/s10115-014-0808-1
47. Olmezogullari, E., Ari, I.: Online Association Rule Mining over Fast Data. In: *Proceedings of the IEEE International Congress on Big Data*, pp. 110–117 (2013). DOI 10.1109/BigData.Congress.2013.77
48. Prasad, B., Agarwal, S.: Stream data mining: Platforms, algorithms, performance evaluators and research trends. *Int. Journal of Database Theory and Application* **9**(9), 201–218 (2016)
49. Ramírez-Gallego, S., García, S., Herrera, F.: Online entropy-based discretization for data streaming classification. *Future Generation Computer Systems* **86**, 59–70 (2018). DOI 10.1016/j.future.2018.03.008

50. Ramírez-Gallego, S., Krawczyk, B., García, S., Woźniak, M., Herrera, F.: A survey on data preprocessing for data stream mining: Current status and future directions. *Neurocomputing* **239**, 39–57 (2017). DOI 10.1016/j.neucom.2017.01.078
51. Saggi, M.K., Jain, S.: A survey towards an integration of big data analytics to big insights for value-creation. *Information Processing & Management* **54**(5), 758–790 (2018). DOI 10.1016/j.ipm.2018.01.010
52. Sobolewski, P., Woźniak, M.: SCR: simulated concept recurrence – a non-supervised tool for dealing with shifting concept. *Expert Systems* **34**(5), e12,059 (2017). DOI 10.1111/exsy.12059
53. Souza, V.M., dos Reis, D.M., Maletzke, A.G., Batista, G.E.: Challenges in benchmarking stream learning algorithms with real-world data. *Data Min Knowl Disc* **34**, 1805–1858 (2020). DOI 10.1007/s10618-020-00698-5
54. Sun, A.Y., Zhong, Z., Jeong, H., Yang, Q.: Building complex event processing capability for intelligent environmental monitoring. *Environmental Modelling & Software* **116**, 1–6 (2019). DOI 10.1016/j.envsoft.2019.02.015
55. Tidke, B., Mehta, R.G., Dhanani, J.: Real-Time Bigdata Analytics: A Stream Data Mining Approach. In: *Proc. 5th International Conference on Recent Findings in Intelligent Computing Techniques*, pp. 345–351 (2018). DOI 10.1007/978-981-10-8636-6_36
56. Uysal, A.K., Gunal, S.: The impact of preprocessing on text classification. *Information Processing & Management* **50**(1), 104–112 (2014). DOI 10.1016/j.ipm.2013.08.006
57. Wang, Y., Gao, H., Chen, G.: Predictive complex event processing based on evolving Bayesian networks. *Pattern Recognition Letters* **105**, 207–216 (2018). DOI 10.1016/j.patrec.2017.05.008
58. Webb, G.I.: Contrary to Popular Belief Incremental Discretization can be Sound, Computationally Efficient and Extremely Useful for Streaming Data. In: *Proc. IEEE International Conference on Data Mining*, pp. 1031–1036 (2014). DOI 10.1109/ICDM.2014.123
59. Webb, G.I., Lee, L.K., Goethals, B., Petitjean, F.: Analyzing concept drift and shift from sample data. *Data Mining and Knowledge Discovery* **32**(5), 1179–1199 (2018). DOI 10.1007/s10618-018-0554-1
60. Zhang, L., Lin, J., Karim, R.: Sliding window-based fault detection from high-dimensional data streams. *IEEE Trans. Syst. Man Cybern. Syst.* **47**(2), 289–303 (2017). DOI 10.1109/TSMC.2016.2585566
61. Zhang, S., Zhang, C., Yang, Q.: Data preparation for data mining. *Applied Artificial Intelligence* **17**(5-6), 375–381 (2003). DOI 10.1080/713827180