



HAL
open science

Machine Learning for a Context Mining Facility

Nourhène Ben Rabah, Manuele Kirsch Pinheiro, Bénédicte Le Grand, Ali Jaffal, Carine Souveyet

► **To cite this version:**

Nourhène Ben Rabah, Manuele Kirsch Pinheiro, Bénédicte Le Grand, Ali Jaffal, Carine Souveyet. Machine Learning for a Context Mining Facility. 16th Workshop on Context and Activity Modeling and Recognition, Mar 2020, Austin, United States. pp.678. hal-02899048

HAL Id: hal-02899048

<https://hal.science/hal-02899048v1>

Submitted on 14 Jul 2020

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Machine Learning for a Context Mining Facility

Nourhène BEN RABAH, Manuele KIRSCH PINHEIRO, Bénédicte LE GRAND, Ali JAFFAL, Carine SOUVEYET
Centre de Recherche en Informatique, Université Paris 1 Panthéon Sorbonne
Paris, France

{ Nourhene.Ben-Rabah | Manuele.Kirsch-Pinheiro | Benedicte.Le-Grand | Ali.Jaffal | Carine.Souveyet }@univ-paris1.fr

Abstract—This paper considers generalizing context reasoning capabilities through a context mining facility offered to all Information System applications. This facility requires mining context data at the system scale, which raises several challenges for Machine Learning approaches used for such mining. Through a detailed literature review, we analyze these approaches with regard to the requirements of such a context mining facility at the Information System level, pointing to the potential and to the challenges raised by this perspective.

Keywords—context mining, context data, machine learning

I. INTRODUCTION

Observing and gathering context information from the physical environment is no longer a challenge. The development of the Internet of Things (IoT) technology and smart devices makes it now possible to easily collect multiple data about people from the applications they use and from their surrounding environments e.g. through sensors. In Information Systems (IS), these data can be exploited for decision making as well as for adaptation and automation purposes. An IS is not a simple set of applications. It is also a set of data, resources (both technical and human), and processes of a given organization. Those work together in order to help this organization fulfill its strategic goals and needs. An IS can be seen as a living organism that must evolve with the organization, whose goals and needs are constantly evolving too. Context data can contribute to the evolution of this ecosystem of applications and resources. For instance, we can easily imagine using context data to adapt business processes, by allocating tasks –or skipping some of them– according to actors’ current context; such data can also help build new Key Performance Indicators for smart cities or digital twin scenarios.

Facing the potential offered by context data, Information Systems might now consider context support as a facility, i.e. as a service offered by the system itself to its applications. Viewing context provision as a facility represents an important shift from the traditional context support paradigm: instead of considering context data for particular applications with precise goals, it becomes necessary to gather as much data as possible and provide appropriate reasoning mechanisms for current but also future applications (i.e. consumers of the context facility).

This new vision represents an exciting challenge for context support. Besides obvious issues related to storage, privacy, security and legal issues, the massive availability of context data also raises some questions regarding Machine Learning (ML) approaches. Indeed, the added value of such context data relies on the capability of extracting knowledge from it, by mining such data for

applications, organizations and individuals. Some features of context data, when considered at an Information System scale, may represent a challenge for ML techniques. First of all, context data is uncertain and heterogeneous by nature. Missing and potentially erroneous data from several different formats should be expected. Quality of data cannot be assured, since it depends on unpredictable environment events (e.g. device or network failure, low battery, human intervention, etc.). Besides, context data is evolutive: new sensors and context sources, with different data types and formats, may integrate the system, while others may disappear. Under such conditions, traditional ML tasks, such as selecting features, identifying relevant testing and training sets, labelling data or just training, may become complex.

Thinking of context as a facility implies collecting data not only for a specific application, but for several potential consumers, including future ones. Therefore, different ML approaches may be considered to address the various uses of such a context facility. For example, extracting indicators (KPI) and tendencies for stakeholders in business organizations will not necessarily be achieved the same way as providing context reasoning to applications for runtime adaptation. Indeed, different ML approaches are required and it is important to understand their strengths and drawbacks in view of these evolutive, uncertain and heterogeneous context data at a large scale.

In this paper, we study the opportunity of integrating ML approaches as part of a context mining facility, transforming context support into a service offered by Information Systems to their applications and users.

The remaining of this paper is organized as follows: Section II presents challenges of the proposed context mining facility, while Section III discusses the use of ML approaches on context data. In Section IV, we discuss the applicability of such approaches to this IS-level facility vision, before concluding in Section V.

II. TOWARDS CONTEXT MINING AS FACILITY

Numerous context-aware applications in the literature use ML techniques for reasoning or “mining” context data and for extracting meaningful information from these data. Examples include applications in health care [1], smart cities [2] and IoT [3] [4], just to cite a few. Most of them focus on specific applications, and demonstrate the interest of applying ML approaches to context data. For example, reasoning mechanisms may allow applications to detect or anticipate particular situations and to adapt their behavior accordingly [5] [6] [7].

Moreover, with the growing development of IoT and sensing technologies, gathering context data becomes

easier, and the potential of mining such data only starts to be foreseen, particularly for Information Systems. In such systems, mining context data may have multiple uses: adaptation of system behavior and applications, recommendation of actions, data prediction, anticipation of user’s needs, and decision making (e.g. [8] [9]). With the availability of context data, we may expect a generalization of such “smart” behaviors, built upon mining solutions. According to [10], we can already observe a trend toward increasingly sophisticated systems, coined as “intelligent”, “context-aware”, “adaptive”, “situated”, etc., for which the notion of ‘context’ is central: they are aware of the context that they are used in, and intelligently adapt to this context at runtime.

However, generalizing such smart behavior implies generalizing context mining to the whole Information System. Making context data available only for the applications composing such system is not enough, since the cornerstone of this expected smart behavior remains the capability of reasoning (and mining) such data. Instead of collecting and mining context data application by application, we consider integrating these tasks to the Information System as a facility, i.e., as a service offered by the system to whatever application that needs it.

Considering context mining as an Information System facility implies an important shift in the way context support is currently handled, since the scale changes drastically. We no longer consider single applications, with precise goals and data, but offer a service for a full ecosystem of users (stakeholders, employees, customers, etc.) and applications supporting different business processes. When considering a given application, its developers/designers may define precisely what context data will be considered, the processing steps these data need and the most adapted ML technique to apply considering these data and the application’s purposes. If we consider a service at the Information System scale, we can neither focus on a single purpose, nor predefine the set of context data to be used, since the same service is supposed to be available to many different applications. It is supposed to satisfy the needs of existing applications, considering currently available sources of context data, but also consider future sources of data and the needs of future applications and users.

Previous works on context-aware computing have considered the evolution of context sources through middleware and context models allowing new sources, data types and formats to be easily added [11] [12] [13]. However, the same does not necessarily apply to reasoning mechanisms. Indeed, as we will discuss in Section III, ML approaches usually consider specific data, which are formatted and prepared specifically for the algorithms that will be applied. Those algorithms may also vary according to the purpose of the analysis (classification, regression, clustering, etc.).

Thus, considering context mining as a facility involves overcoming such specificities in order to propose a general service, which must also cope with context data characteristics. Context-aware computing literature [11] [14] highlights multiple features of context data,

summarized in Table I, which make such data particularly challenging for some ML approaches. Context data is naturally uncertain and incomplete; it may contain errors and be very dynamic; it is heterogeneous, including different formats (numeric and symbolic, structured and unstructured, etc.), types and sources; it may be observed using different frequencies or be pushed up by its sources (i.e. sensors). The sources of these data may also vary, as new sources can be integrated into the system, while others may disappear (temporarily or definitely). All this suggests new data and potentially new data formats for ML algorithms, which may be unable to handle them.

TABLE I. CONTEXT DATA MAIN FEATURES

| Characteristic | Definition |
|----------------|---|
| Uncertainty | Context data may contain errors and imprecisions. Data quality cannot be assured. |
| Incompleteness | Context data may be incomplete; we cannot guarantee that 100% of possible observations have been recorded. |
| Heterogeneity | Context data may be gathered from multiple sources, using multiple formats even for the same data. |
| Dynamicity | Context data may evolve quickly; new data may arrive frequently, and data may be observed with different frequencies. |

Therefore, considering context mining as a facility implies requirements that may impact the reasoning techniques at stake and notably those based on ML approaches (see Table II).

TABLE II. REQUIREMENTS FOR A CONTEXT MINING FACILITY

| Requirement | Definition |
|-------------|--|
| R0 | Guaranteeing security and privacy of context data. |
| R1 | Supporting context data features (Table I). |
| R2 | Supporting multiple heterogeneous context data sources. |
| R3 | Supporting the evolution of context data sources and formats; new unexpected or non-predefined context sources and formats should be easily integrated in the system. |
| R4 | Supporting dynamicity of context data sources; these sources may become offline unexpectedly, for small or large periods of time, they may also totally disappear, and other new sources from known and unknown formats and data types may integrate the system. |
| R5 | Supporting online processing; the Information System should constantly remain running, and since critical applications may depend on context data, high availability is mandatory for supporting system applications. |
| R6 | Operating with no or little human intervention; considering frequent human intervention for keeping the system running, for cleaning or preparing new context data is unfeasible considering the volume of data, the dynamicity of the data set and the need for availability. |

The requirements listed in Table II consider a heterogeneous and constantly evolving data set formed by observed context data. Under such conditions, it is difficult to assume a previous knowledge about these data. Moreover, stopping the system for data preparation or preprocessing pre-treatment or data preparing may be

costly in terms of consequences for the system applications that use the service.

It is thus necessary to examine the impact of these requirements on ML approaches in order to achieve the smart behavior promoted by context-awareness at an Information System scale. In order to tackle this issue, we analyzed how ML approaches are currently used for mining context data and what their requirements are for correctly analyzing these data. Our goal is to evaluate whether our view of context mining facility is possible and which challenges should be tackled for applying ML approaches to context data at this scale.

III. USING MACHINE LEARNING FOR CONTEXT MINING

As explained previously, the vision we propose of context as an Information System facility requires the collaboration of several research areas, among which context-aware computing and ML. We have therefore considered various research communities in our literature review. We have first studied the articles published at the CoMoRea workshop, as it is dedicated to context modeling and reasoning. The papers we have selected from this conference reflect very interesting contributions. We have however noticed that many of them are application-specific (and therefore context-specific), whereas we seek more general solutions. Moreover, the justification of choice of the algorithm used for the reasoning part is not always very detailed in the papers we studied.

We have thus widened the spectrum of our literature review in order to include the ML and IoT communities. Instead of focusing on specific conferences, we have performed keywords-based queries like “context prediction” and “context awareness”, with a precision on the research area, like “ML”, “data mining”, or “IoT”. We have discarded papers published before 2013 as we wanted to focus on recent works. These keyword-based queries returned more than 200 research papers. We read their abstracts and selected those which went further than merely “provide context data” but also included some reasoning or mining mechanisms. Among the remaining articles, we focused on the approaches that seemed appropriate in terms of scalability and discarded some rule-based or ontology-based solutions, focusing on works using ML approaches only.

Indeed, since our main goal here is to better understand challenges of ML for a context facility use, we have focused only on these approaches. We are aware that ontology and ML can be combined. Several papers propose to facilitate ontology generation thanks to ML approaches [15] [16], while a few others try to use ontologies as knowledge representations to support ML approaches [17]. However, they still face some challenges (that we underline later in this section), in addition to some extra issues such as ontology learning, or temporal reasoning, as pointed by [16]. Thus, in order to focus on ML issues, we decided to focus our analysis on ML approaches, ignoring those that combine them with other reasoning techniques.

Following this methodology, we finally selected 30 papers that we analyzed in depth. These papers have been

published in various communities, covering a large spectrum of expertise. About 20% come from the context modeling and reasoning community; 23% from sensors and mobile networks area; 20% have been published in ML journals of conferences. Finally, 20% come from generalist computer science sources and about 10% from areas outside the computer science fields, such as biomedicine or social sciences. As mentioned earlier, all papers have been published since 2013; half of them have even been published since 2016.

Whatever the adopted paradigm (supervised, unsupervised, semi-supervised or reinforcement learning [18] [19]), ML is about selecting an algorithm and training it on some data. The effectiveness of a given method depends on various factors, such as the quality of the training data, the chosen algorithm and its hyperparameters. Low quality data may compromise the success of the most powerful ML algorithms [20] [21]. As mentioned earlier, raw data obtained from heterogeneous sensors may be noisy, scattered and even incomplete. This may lead to various difficulties such as increased processing time, higher model complexity and overfitting.

In order to address these problems, many studies in context recognition include a pre-processing step. For example, in [7], the authors proposed a general context prediction structure based on a pre-processing phase and a context prediction phase. In [18], the authors presented a framework based on four steps: data processing, feature set generation, model selection and model combination for predicting the consumption of air conditioning in residential buildings. We detail frequently-used pre-processing tasks below.

Data cleaning techniques use one or more filters to identify noisy data and to correct or delete them [22] [23] [24]. They can also process missing values and detect outliers [18] [25].

Data transformation methods modify data representation to make them suitable as model inputs. They involve digitalization and normalization: digitalization encodes qualitative data into numerical data. Indeed, some algorithms can work directly with qualitative data, such as K-Nearest Neighbors (KNN), Naive Bayes (NB), Decision Tree (DT), or Random Forest (RF). However, many others require numerical input and output variables to operate properly. Normalization modifies the values of numerical data into a common scale. In deep neural networks, if numerical data do not have similar ranges of values then they will have a negative influence on gradient descent optimization methods, with a lower learning rate [25] [26].

Feature extraction and selection techniques identify relevant information and help remove as much irrelevant and redundant features as possible from raw data. If there are not enough informative features, the model will not be able to accomplish its task. If there are too many features or irrelevant ones, then the model will be more resource consuming and harder to train; practitioners agree that most of the time spent in building a ML pipeline is dedicated to feature engineering [27].

Data augmentation methods create additional training samples since some ML algorithms such as deep neural

networks need a huge amount of data to learn effectively [28]. However, collecting such training data is often expensive and laborious.

Unbalanced data processing methods are used to address the problem of class imbalance (i.e. when there is a disproportionate ratio of instances in each class) [29]. This is often the case with real world data, and the models learned from them generally have a good accuracy on the majority class but perform poorly on other classes.

These various pre-processing techniques have been used in several research studies for context modelling and recognition. In Table III, we compare the solutions proposed in the papers of our literature review according to the data processing mechanisms employed by each work, as well as the ML approach that has been used. We pay special attention to the following pre-processing tasks: cleaning (*c1*), transformation (*c2*), feature extraction and selection (*c3*), data augmentation (*c4*), and unbalanced data processing (*c5*). Regarding the ML approach, we indicate the algorithm(s) (*c7*) that have been used. We also report whether or not the authors mentioned hyperparameters tuning (*c6*). A hyperparameter is a parameter of the ML algorithm and not of the model. Setting hyperparameters can improve the performance of algorithms. Only 67 % of the 30 selected papers appear in Table III as we discarded the survey papers, which did not describe original experimentations.

TABLE III. COMPARATIVE OVERVIEW OF CONTEXT RECOGNITION APPROACHES BASED ON MACHINE LEARNING

| Ref. | Data criteria | | | | | Reasoning criteria | |
|------|---------------|-----------|-----------|-----------|-----------|--------------------|--|
| | <i>c1</i> | <i>c2</i> | <i>c3</i> | <i>c4</i> | <i>c5</i> | <i>c6</i> | <i>c7</i> |
| [23] | ✓ | ✗ | ✓ | ✗ | ✗ | ✓ | SVM |
| [30] | ✗ | ✗ | ✓ | ✗ | ✗ | ✗ | KNN, Gaussian Naive Bayes, DT, RF, RMD |
| [31] | ✗ | ✓ | ✓ | ✗ | ✗ | ✗ | NB, DT, RF, SVM, KNN, Adaptive Boosting, LR, ANN |
| [32] | ✓ | ✓ | ✓ | ✗ | ✗ | ✓ | MLP |
| [33] | ✓ | ✓ | ✓ | ✗ | ✗ | ✗ | NB, ANN, Bayesian Network |
| [34] | ✗ | ✗ | ✓ | ✗ | ✗ | ✓ | KNN |
| [35] | ✗ | ✗ | ✓ | ✗ | ✗ | ✗ | KNN, NB, RF |
| [36] | ✓ | ✗ | ✓ | ✗ | ✓ | ✗ | NB, DT, LR, Adaboost, SVM |
| [18] | ✓ | ✗ | ✓ | ✗ | ✗ | ✓ | Support Vector Regression, Ensemble tree, ANN |
| [37] | ✗ | ✗ | ✓ | ✗ | ✗ | ✗ | RF, SVM, NB, Random Tree, Bayesian Network |
| [38] | ✗ | ✗ | ✓ | ✗ | ✗ | ✗ | DT, MLP, LR |
| [39] | ✓ | ✗ | ✓ | ✗ | ✗ | ✗ | MLP, SVM, LogitBoost |
| [40] | ✗ | ✗ | ✓ | ✗ | ✗ | ✗ | RF |
| [22] | ✓ | ✗ | ✓ | ✗ | ✗ | ✓ | RF |
| [24] | ✓ | ✗ | ✗ | ✗ | ✗ | ✓ | CNN+ symbolic model |
| [28] | ✗ | ✓ | ✗ | ✓ | ✗ | ✓ | CNN |
| [26] | ✓ | ✓ | ✗ | ✗ | ✗ | ✓ | CNN+ LSTM |
| [25] | ✗ | ✗ | ✗ | ✓ | ✓ | ✓ | CNN |
| [29] | ✗ | ✗ | ✗ | ✗ | ✓ | ✗ | RNN |
| [41] | ✗ | ✗ | ✗ | ✗ | ✗ | ✓ | CNN |

(*) No (✓) Yes

We note that the authors of [23] propose a multi-class Support Vector Machine (SVM) based on features extracted from both an accelerometer and a gyroscope after a pre-processing step for noise reduction with a median filter and Butterworth filter. In [30], the authors propose a multi-class classifier to determine the position of smartphones in different contexts of use. They apply KNN, Gaussian Naive Bayes, DT, RF and Random Mixture Model (RMD), and report that KNN provide better performance than the other algorithms. For data pre-processing, the approach is based on a windowing step and a feature extraction step. In [32], the authors suggest a framework based on genetic algorithm to extract relevant features from IoT raw data and a Multilayer Perceptron (MLP) to classify these data in smart industrial applications.

To obtain better performance, several studies use ensemble methods such as [18] that combines the predictions of three models resulting from Support Vector Regression, Ensemble Tree and Artificial Neural Network (ANN) algorithms for predicting the consumption of air conditioning in residential buildings. To ensure that their data meet the input requirements of the models, they perform a linear interpolation for missing and incorrect data. Then, to select the right feature set, they use a statistical measure. In [37], the authors propose an approach for detecting the current transportation mode of a user from his/her smartphone sensors data. They propose to divide the collected data into consecutive non-overlapping time sequences and to extract four features for each sequence and each sensor. Then, they combine multiple learners to improve their performance. In [38], the authors present an ensemble method that combines the predictions of three models resulting from DT, MLP, and Logistic Regression (LR) for human activity recognition. To determine the class of a new activity, they consider the predictions (i.e. classes) of the three models, and choose the class with the highest number of votes. Their results show that ensemble learning can achieve significant improvements for activity recognition when compared to what each learning algorithm can achieve individually. The same problem is also investigated in [39]; in this case, however, the authors combine the results of other classifiers such as MLP, SVM, and LogitBoost. In addition, they use a clustering method to select 18 relevant features from 24 features and they obtain a good accuracy of 91.15%. In [22], the authors introduce a multi-class classification approach based on ultra-wide band sensor measurements and RF to detect when old people fall down. The pre-processing phase includes filtering, feature extraction, stream windowing, change detection and buffering. The classifier obtains the lowest error rate by setting the number of trees at 200.

To extract relevant features without effort and improve their performance, several studies use different deep neural networks architectures. The authors in [25] [28] [41] [42] [43] propose a Convolutional Neural Network (CNN) allowing feature extraction and classification for human activity recognition. The same problem is also investigated in [26], where the authors propose a generic deep framework for activity recognition based on convolutional and Long Short-Term Memory (LSTM) recurrent units. In

[34], the authors present a learning deep features for KNN to improve the classification performance. In these works, it is not necessary to extract the hand-crafted features or to use statistical methods or frequency transformation coefficients [34], as deep features can be extracted using deep learning approaches. The first layers of networks extract features that the following layers will combine to form increasingly complex and abstract concepts.

To ensure the performance of deep learning networks, the authors of [26] pre-process sensor data to fill in the missing values by linear interpolation and to perform channel normalisation to the interval [0,1]. For [25], a normalization step is required for the raw signal extracted from the accelerometers to have a common scale. They propose to apply a mean-zero normalisation. In [28], the authors use data augmentation methods such as Gaussian noise to artificially create new training data from existing learning data. In [29], the authors present a Recurrent Neural Network (RNN) based on a windowing approach for human activity recognition. They apply a synthetic minority over sampling technique to deal with the class imbalance problem.

In the next section, we analyze lessons learned from these ML approaches face to the challenges proposed by a context mining facility.

IV. MACHINE LEARNING APPROACHES FACE TO CONTEXT MINING CHALLENGES

The comparative table we proposed above (Table III), based on our literature review, allows us to make several observations. We observe that the authors from context and from sensors/mobile networks communities are up to date with regard to the state of the art in ML solutions. Indeed, there is no striking difference in the algorithms used in those communities as compared to the ML one. However, we notice that most of these “non machine learning experts” used the algorithms without mentioning any hyperparameters tuning phase. This may suggest that they could benefit from experts help in this area to better use existing ML solutions, notably considering hyperparameters optimization. Besides, we can note that deep learning solutions (among which CNN, frequently seen in Table III) are becoming more and more popular. Indeed, one of their strengths is that they allow skipping the feature extraction and selection step. Moreover, these approaches can handle large volumes of data, which is one of the requirements for context data.

Another observation we can make is that very few works consider the variety of context data and take into account all their features. The approaches we have reviewed apply ML techniques to subsets of very specific data, such as a predefined set of sensors.

As we may observe in Table III, ML approaches illustrated by our literature review heavily rely on data pre-processing phases. The quality of these approaches depends on these phases, which may be mandatory in some cases. Focusing on precise context data allows the execution of these pre-processing phases, since data types and formats are known in advance. However, when considering the context mining facility requirements highlighted in Table II, we may note that the execution of such phases cannot be

guaranteed. Since new context data sources and formats may join the system at any moment, these pre-processing phases can be put in question. On the one side, if these phases are not reconsidered, new relevant data may remain ignored and context data unexplored. On the other side, stopping the facility for re-executing those may also have negative consequences on critical applications that may depend on it. This interruption does not concern only the training phase, but all the preprocessing tasks mentioned in Section III, and notably data augmentation and unbalanced data processing.

TABLE IV. ESTIMATION OF THE IMPACT OF CONTEXT MINING FACILITY REQUIREMENTS ON ML DATA CRITERIA

| ML Data criteria | Context mining facility requirements | | | | | |
|-----------------------|--------------------------------------|----|----|----|----|----|
| | R1 | R2 | R3 | R4 | R5 | R6 |
| Cleaning | ⊖ | ⊖ | | ⊖ | ⊖ | ⊖ |
| Transformation | ⊖ | ⊖ | ⊖ | | ⊖ | |
| Feature extraction | | ⊖ | ⊖ | ⊖ | ⊖ | ⊖ |
| Data augmentation | ⊖ | | ⊖ | ⊖ | ⊖ | |
| Proc. unbalanced data | ⊖ | | ⊖ | ⊖ | | |

(⊖) Negative impact estimated

Table IV confronts context mining facility requirements summarized in Table II, and ML practices reported in Table III (since we do not analyze security and privacy aspects here, requirement R0 is not considered in Table IV). We may observe that requirements in Table II make some steps preconized by most of the approaches previously discussed more difficult to achieve. For instance, processing unbalanced data can be challenging when considering context data characteristics (R1), and notably uncertainty and incompleteness. Similarly, feature extraction and selection can be complex without human intervention or previous knowledge about the data. To sum up, Table IV highlights the fact that, although ML algorithms have proven to be useful for mining specific context data, the overall process necessary for applying those can be challenging when considering a large scale. This is true not only for supervised approaches, but also for semi-supervised ones, since they also require human intervention somewhere in the pipeline. Considering the IS scale, even a minimal human intervention may be complex. Similarly, unsupervised approaches are impacted too since they also rely on pre-processing. Indeed, any ML algorithm will fail to discover a hidden pattern or trend from raw data if that data is inadequate, irrelevant or incomplete.

V. CONCLUSIONS

Recent advances in middleware solutions allow to integrate new context data sources at runtime [44] [45]. The availability of such data raises multiple issues: how can ML algorithm exploit these data? How can we make such reasoning capabilities available to whatever application in an Information System? How can we make ML approaches scale up to a system (and not a precise application) scale?

In this paper, we presented a literature review of works that applied ML techniques to context data. Through this study, we could observe that several approaches were not general enough and often focused on specific types of

context data, ignoring the heterogeneity of such data. Most of these works consider context data as a “traditional” data, and do not fully take into account their features. Besides, the scale involved in a context mining facility also implies additional requirements, since it opens the possibility of an evolving set of applications, customers and context data, rather than precise and well-identified applications.

These observations highlight the challenges of applying traditional ML process in the case of a context mining facility. However, the use of ML for reasoning in this context is not impossible. Although ML techniques are currently evolving, they still require a sequence of pre-processing tasks that are hardly scalable. Their performance is also very sensitive to some design decisions such as algorithm selection, hyperparameters tuning, etc. Therefore, in order to reach a true context mining facility, we should evolve these practices towards powerful tools that gradually remove the human from the loop. We can already observe this tendency through initiatives such as automated ML (AutoML) frameworks [46], which aim at automating the entire ML pipeline. Among these frameworks, we can mention commercial solutions, hosted by leading cloud providers such as Amazon Machine Learning, Microsoft Azure Machine Learning and Google Cloud AutoML, as well as academic and open source frameworks such as Auto-WEKA [47] or autosklearn [48]. However, the latter only allow hyperparameters tuning and algorithms selection. We strongly believe that the future of a large scale context mining lies in the automatization of preprocessing phases and in the development of configurable frameworks that remain parametrizable according to application and organization needs. The potential of generalizing such reasoning capabilities is huge on different application areas, such as Smart cities and Industry 4.0. We are convinced that a stronger collaboration between ML experts and context specialists could help make ML solutions more flexible and adapted to context data, and further help reaching the full potential of context data for large Information Systems.

REFERENCES

- [1] B. Yuan and J. Herbert, "Context-aware Hybrid Reasoning Framework for Pervasive Healthcare," *Personal Ubiquitous Comput.*, vol. 18, no. 4, pp. 865-881, 2014.
- [2] A. K. Ramakrishnan, D. Preuveneers and Y. Berbers, "A Bayesian Framework for Life-Long Learning in Context-Aware Mobile Applications," in *Context in Computing: A Cross-Disciplinary Approach for Modeling the Real World*, P. Brézillon and A. J. Gonzalez, Eds., Springer NY, 2014, pp. 127-141.
- [3] A. M. Otebolaku and G. M. Lee, "Towards context classification and reasoning in IoT," in *14th International Conference on Telecommunications (ConTEL)*, 2017.
- [4] H. Rahman, R. Rahmani and T. Kanter, "Multi-Modal Context-Aware reasoner (CAN) at the Edge of IoT," *Procedia Computer Science*, vol. 109, pp. 335 - 342, 2017.
- [5] R. Mayrhofer, "Context Prediction based on Context Histories: Expected Benefits, Issues and Current State-of-the-Art," in *1st International Workshop on Exploiting Context Histories in Smart Environments (ECHISE 2005)*, *3rd International Conference on Pervasive Computing (PERVASIVE 2005)*, 2005.
- [6] S. Sigg, S. Haseloff and K. David, "An alignment approach for context prediction tasks in ubicomp environments," *IEEE Pervasive Computing*, vol. 9, no. 4, pp. 90-97, 2010.
- [7] D. Ameyed, M. Miraoui and C. Tadj, "A survey of prediction approach in pervasive computing," *International Journal of Scientific & Engineering Research*, vol. 6, pp. 306-316, 2015.
- [8] S. Najar, M. Kirsch-Pinheiro and C. Souveyet, "Service discovery and prediction on Pervasive Information System," *J. of Ambient Intelligence and Humanized Comp.*, vol. 6, no. 4, pp. 407-423, 2015.
- [9] C. D. Maio, G. Fenza, V. Loia, F. Orciuoli and E. Herrera-Viedma, "A framework for context-aware heterogeneous group decision making in business processes," *Knowledge-Based Systems*, vol. 102, pp. 39 - 50, 2016.
- [10] C. Bauer and A. K. Dey, "Considering context in the design of intelligent systems: Current practices and suggestions for improvement," *Journal of Systems and Software*, vol. 112, pp. 26-47, 2016.
- [11] C. Bettini, O. Brdiczka, K. Henriksen, J. Indulska, D. Nicklas, A. Ranganathan and D. Riboni, "A survey of context modelling and reasoning techniques," *Pervasive and Mobile Computing*, vol. 6, no. 2, pp. 161-180, Apr 2010.
- [12] N. Paspallis and G. A. Papadopoulos, "A Pluggable Middleware Architecture for Developing Context-aware Mobile Applications," *Personal Ubiquitous Comp.*, vol. 18, no. 5, pp. 1099-1116, 2014.
- [13] M. Wagner, R. Reichle and K. Geihs, "Context as a service - Requirements, design and middleware support," in *Pervasive Computing and Communications Workshops (PERCOM Workshops)*, *2011 IEEE International Conference on*, 2011.
- [14] K. Henriksen, J. Indulska and A. Rakotonirainy, "Modeling context information in pervasive computing systems," in *LNCS 2414 - First International Conference in Pervasive Computing (Pervasive 2002)*, 2002.
- [15] D. Riboni and C. Bettini, "COSAR: hybrid reasoning for context-aware activity recognition," *Personal and Ubiquitous Computing*, vol. 15, pp. 271-289, 2011.
- [16] G. Civitarese, R. Presotto and C. Bettini, "Context-driven Active and Incremental Activity Recognition," *arXiv preprint arXiv:1906.03033*, 2019.
- [17] M. A. Razzaq, M. B. Amin and S. Lee, "An ontology-based hybrid approach for accurate context reasoning," in *2017 19th Asia-Pacific Network Operations and Management Symposium (APNOMS)*, 2017.
- [18] C. Lork, B. Rajasekhar, C. Yuen and N. M. Pindoriya, "How many watts: A data driven approach to aggregated residential air-conditioning load forecasting," in *CoMoRea 2017, IEEE International Conference on Pervasive Computing and Communications Workshops (PerCom Workshops)*, 2017.
- [19] O. B. Sezer, E. Dogdu and A. M. Ozbayoglu, "Context-aware computing, learning, and big data in internet of things: a survey," *IEEE Internet of Things Journal*, vol. 5, pp. 1-27, 2017.
- [20] S. García, J. Luengo and F. Herrera, "Tutorial on practical tips of the most influential data preprocessing algorithms in data mining," *Knowledge-Based Systems*, vol. 98, pp. 1-29, 2016.
- [21] S. Ramírez-Gallego, B. Krawczyk, S. García, M. Woźniak and F. Herrera, "A survey on data preprocessing for data stream mining: Current status and future directions," *Neurocomputing*, vol. 239, pp. 39-57, 2017.
- [22] G. Mokhtari, Q. Zhang and A. Fazlollahi, "Non-wearable UWB sensor to detect falls in smart home environment," in *CoMoRea 2017, IEEE International Conference on Pervasive Computing and Communications Workshops (PerCom Workshops)*, 2017.
- [23] D. Anguita, A. Ghio, L. Oneto, X. Parra Perez and J. L. Reyes Ortiz, "A public domain dataset for human activity recognition using smartphones," in *Proceedings of the 21th International European Symposium on Artificial Neural Networks, Computational Intelligence and Machine Learning*, 2013.
- [24] F. M. Rueda, S. Lüdtke, M. Schröder, K. Yordanova, T. Kirste and G. A. Fink, "Combining Symbolic Reasoning and Deep Learning for Human Activity Recognition," in *CoMoRea 2019, IEEE International Conference on Pervasive Computing and Communications Workshops (PerCom Workshops)*, 2019.
- [25] M. Zeng, L. T. Nguyen, B. Yu, O. J. Mengshoel, J. Zhu, P. Wu and J. Zhang, "Convolutional neural networks for human activity

- recognition using mobile sensors," in *6th Int. Conf. on Mobile Computing, Applications and Services*, 2014.
- [26] F. Ordóñez and D. Roggen, "Deep convolutional and lstm recurrent neural networks for multimodal wearable activity recognition," *Sensors*, vol. 16, p. 115, 2016.
- [27] A. Zheng and A. Casari, *Feature engineering for machine learning: principles and techniques for data scientists*, "O'Reilly Media, Inc.", 2018.
- [28] R. Grzeszick, J. M. Lenk, F. M. Rueda, G. A. Fink, S. Feldhorst and M. Hompel, "Deep neural network based human activity recognition for the order picking process," in *4th Int. Workshop on Sensor-based Activity Recognition and Interaction*, 2017.
- [29] F. Al Machot, S. Ranasinghe, J. Plattner and N. Jnoub, "Human activity recognition based on real life scenarios," in *CoMoRea 2018, IEEE International Conference on Pervasive Computing and Communications Workshops (PerCom Workshops)*, 2018.
- [30] I. Diaconita, A. Reinhardt, F. Englert, D. Christin and R. Steinmetz, "Do you hear what i hear? using acoustic probing to detect smartphone locations," in *CoMoRea 2014, IEEE International Conference on Pervasive Computing and Communication Workshops (PERCOM WORKSHOPS)*, 2014.
- [31] I. H. Sarker, A. S. M. Kayes and P. Watters, "Effectiveness analysis of machine learning classification models for predicting personalized context-aware smartphone usage," *Journal of Big Data*, vol. 6, p. 57, 2019.
- [32] N. Mishra, C.-C. Lin and H.-T. Chang, "A cognitive adopted framework for IoT big-data management and knowledge discovery prospective," *International Journal of Distributed Sensor Networks*, vol. 11, p. 718390, 2015.
- [33] S. Saeedi, A. Moussa and N. El-Sheimy, "Context-aware personal navigation using embedded sensor fusion in smartphones," *Sensors*, vol. 14, pp. 5742-5767, 2014.
- [34] S. Sani, N. Wiratunga and S. Massie, "Learning deep features for kNN-based human activity recognition," in *ICCB*, 2017.
- [35] M. Miettinen, S. Heuser, W. Kronz, A.-R. Sadeghi and N. Asokan, "ConXsense: automated context classification for context-aware access control," in *Proceedings of the 9th ACM symposium on Information, computer and communications security*, 2014.
- [36] L. D. Turner, S. M. Allen and R. M. Whitaker, "Push or delay? decomposing smartphone notification response behaviour," in *Human Behavior Understanding*, Springer, 2015, pp. 69-83.
- [37] L. Bedogni, M. Di Felice and L. Bononi, "Context-aware Android applications through transportation mode detection techniques," *Wireless communications and mobile computing*, vol. 16, pp. 2523-2541, 2016.
- [38] C. Catal, S. Tufekci, E. Pirmir and G. Kocabag, "On the use of ensemble of classifiers for accelerometer-based activity recognition," *Applied Soft Comp.*, vol. 37, pp. 1018-1022, 2015.
- [39] A. Bayat, M. Pomplun and D. A. Tran, "A study on human activity recognition using accelerometer data from smartphones," *Procedia Computer Science*, vol. 34, pp. 450-457, 2014.
- [40] Z. Feng, L. Mo and M. Li, "A Random Forest-based ensemble method for activity recognition," in *2015 37th Annual International Conference of the IEEE Engineering in Medicine and Biology Society (EMBC)*, 2015.
- [41] D. Ravi, C. Wong, B. Lo and G.-Z. Yang, "Deep learning for human activity recognition: A resource efficient implementation on low-power devices," in *IEEE 13th Int. Conf. on Wearable and Implantable Body Sensor Networks (BSN)*, 2016.
- [42] J. Yang, M. N. Nguyen, P. P. San, X. L. Li and S. Krishnaswamy, "Deep convolutional neural networks on multichannel time series for human activity recognition," in *24th Int. Joint Conf. on Artificial Intelligence*, 2015.
- [43] C. A. Ronao and S.-B. Cho, "Deep convolutional neural networks for human activity recognition with smartphone sensors," in *Int. Conf. on Neural Information Processing*, 2015.
- [44] P. Pradeep and S. Krishnamoorthy, "The MOM of context-aware systems: A survey," *Computer Communications*, vol. 137, pp. 44 - 69, March 2019.
- [45] P. Temdee and R. Prasad, *Context-Aware Communication and Computing: Applications for Smart Environment*, Springer, 2018.
- [46] F. Hutter, L. Kotthoff and J. Vanschoren, *Automated Machine Learning-Methods, Systems, Challenges*, Springer, 2019.
- [47] L. Kotthoff, C. Thornton, H. H. Hoos, F. Hutter and K. Leyton-Brown, "Auto-WEKA 2.0: Automatic model selection and hyperparameter optimization in WEKA," *The Journal of Machine Learning Research*, vol. 18, pp. 826-830, 2017.
- [48] M. Feurer, A. Klein, K. Eggenberger, J. Springenberg, M. Blum and F. Hutter, "Efficient and Robust Automated Machine Learning," in *Advances in Neural Information Processing Systems 28*, C. Cortes, N. D. Lawrence, D. D. Lee, M. Sugiyama and R. Garnett, Eds., Curran Associates, Inc., 2015, pp. 2962-2970.