

Implementation of Concept Drifts in Data Mining Process

Ganapam Venkata Suresh Reddy¹ Pasupula Venkata Nagendra Kumar²

¹System Trainee ²Sr.Tech Associate

¹Polaris Financial Technology Ltd. ²Bank of America

Abstract—Eventhough most business processes change over time, contemporary process mining techniques tend to analyze these processes as if they are in a steady state. Processes may change suddenly or gradually. The drift may be periodic (e.g., because of seasonal influences) or one-of-a-kind (e.g. the effects of new legislation). For the process management, it is crucial to discover and understand such concept drifts in processes. This paper presents a generic framework and specific techniques to detect when a process changes and to localize the parts of the process that have changed. Different features are proposed to characterize relationships among activities. These features are used to discover differences between successive populations. The approach has been implemented as a plug-in of the ProM process mining framework and has been evaluated using both simulated event data exhibiting controlled concept drifts and real-life event data from a Dutch municipality.

Keywords—Concept drift, liveness, hypothesis tests, process changes, process mining

I. INTRODUCTION

BUSINESS processes are nothing more than logically related tasks that use the resources of an organization to achieve a defined business outcome. Business processes can be viewed from a number of perspectives, including the control flow, data, and the resource perspectives. In today's dynamic marketplace, it is increasingly necessary for enterprises to streamline their processes so as to reduce cost and to improve performance. In addition, today's customers expect organizations to be flexible and adapt to changing circumstances. New legislations such as the WABO act [1] and the Sarbanes-Oxley Act [2], extreme variations in supply and demand, seasonal effects, natural calamities and disasters, deadline escalations [3], and so on, are also forcing organizations to change their processes. For example, governmental and insurance organizations reduce the fraction of cases being checked when there is too much of work in the pipeline. As another example, in a disaster, hospitals, and banks change their operating procedures. It is evident that the economic success of an organization is more and more dependent on its ability to react and adapt to changes in its operating environment. Therefore, flexibility and change have been studied in-depth in the context of business process management (BPM). For example, process-aware information systems (PAISs) [4] have been extended to be able to flexibly adapt to changes in the process. State-of-the-art workflow management (WFM) and BPM systems [5] provide such flexibility, e.g., we can easily release a new version of a process. In addition, in processes not driven by WFM/BPM systems (such as the usage of medical systems) there is even more flexibility as processes are controlled by people rather than information systems.

Many of today's information systems are recording an abundance of event logs. Process mining is a relatively young research discipline aimed at discovering, monitoring, and improving real processes by extracting knowledge from event logs [6] (Section II-A for a brief introduction). Although flexibility and change have been studied in-depth in the context of WFM and BPM systems, contemporary process mining techniques assume the processes to be in a steady state. For example, when discovering a process model from event logs, it is assumed that the process at the beginning of the recorded period is the same as the process at the end of the recorded period. Using ProM, we have analyzed processes in more than 100 organizations. These practical experiences show that it is very unrealistic to assume that the process being studied is in a steady state. As mentioned earlier, processes may change to adapt to changing circumstances. Concept drift refers to the situation in which the process is changing while being analyzed. There is a need for techniques that deal with such second-order dynamics. Analyzing such changes is of utmost importance when supporting or improving operational processes and to obtain an accurate insight on process executions at any instant of time. When dealing with concept drifts in process mining, the following three main challenges emerge.

1) Change point detection:

The first and most fundamental problem is to detect concept drift in processes, i.e., to detect that a process change has taken place. If so, the next step is to identify the time periods at which changes have taken place. For example, by analyzing an event log from an organization (deploying seasonal processes), we should be able to detect that process changes happen and that the changes happen at the onset of a season.

2) Change localization and characterization:

Once a point of change has been identified, the next step is to characterize the nature of change, and identify the region(s) of change (localization) in a process. Uncovering the nature of change is a challenging problem that involves both the identification of change perspective (e.g., control flow, data, resource, sudden, gradual, and so on) and the identification of the exact change itself. For instance, in the example of a seasonal process, the change could be that more resources are deployed or that special offers are provided during holiday seasons.

3) Change process discovery:

Having identified, localized, and characterized the changes, it is necessary to put all of these in perspective. There is a need for techniques/tools that exploit and relate these discoveries. Unraveling the evolution of a process should result in the discovery of the change process describing these second-order dynamics. For instance, in the example of a seasonal process, we could identify that the process recurs every season. In addition, we can show an animation on how the

process evolved over a period with annotations showing several perspectives such as the performance metrics (service levels, throughput time, and so on) of process at different instances of time.

We can differentiate between two broad classes of dealing with concept drifts when analyzing event logs (Fig. 1).

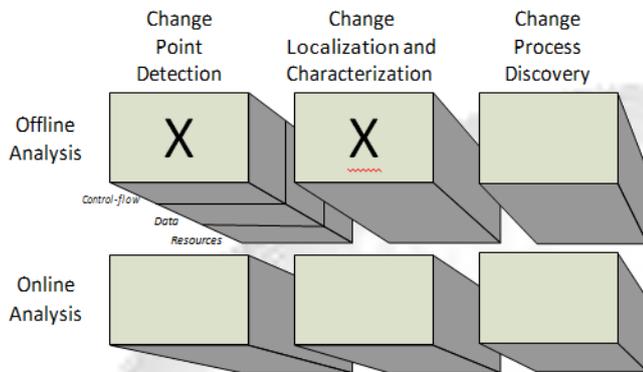


Fig. 1:

Different dimensions of concept drift analysis in process mining.

1) Offline analysis: This refers to the scenario where the presence of changes or the occurrence of drifts need not be uncovered in a real time. This is appropriate in cases where the detection of changes is mostly used in postmortem analysis, the results of which can be considered when designing/improving processes for later deployment. For example, offline concept drift analysis can be used to better deal with seasonal effects (hiring less staff in summer or skipping checks in the weeks before Christmas).

2) Online analysis: This refers to the scenario where changes need to be discovered in near real time. This is appropriate in cases where an organization would be more interested in knowing a change in the behavior of their customers or a change in demand as and when it is happening. Such real-time triggers (alarms) will enable organizations to take quick remedial actions and avoid any repercussions.

In this paper, we focus on two of the challenges: 1) change (point) detection and change localization and 2) characterization in an offline setting (Fig. 1). We define different features and propose a framework for dealing with these two problems from a control-flow perspective. Initially, we show the promise of the techniques proposed in this paper on a synthetic log and later evaluate them on a real-life case study from a large Dutch municipality.

The rest of this paper is organized as follows. Section II provides background on process mining and concept drifts in data mining. Related work is presented in Section III. Section IV describes the various aspects and nature of change, whereas Section V presents the basic idea for change detection in event logs. Section VI introduces various features that capture the characteristics of event logs. Section VII illustrates the significance of statistical hypothesis tests for detecting drifts. Section VIII presents the framework for dealing with concept drifts in process mining, whereas Section IX presents the realization of the proposed approaches in the ProM framework. Section X describes the effectiveness of the

features and the techniques proposed in this paper on a synthetic log as well as a real-life case study. Finally, this paper is summarized with a conclusion and an outlook on some of the open research questions in Section XI.

II. BACKGROUND

In this section, we discuss the basic concepts in process mining and concept drifts in data mining/machine learning.

A. Process Mining

Process mining serves a bridge between data mining and business process modeling [6]. Business processes leave trails in a variety of data sources (e.g., audit trails, databases, and transaction logs). Process mining aims at discovering, monitoring, and improving real processes by extracting knowledge from event logs recorded by a variety of systems (ranging from sensor networks to enterprise information systems). The starting point for process mining is an event log, which is a collection of events. We assume that events can be related to process instances (often called cases) and are described by some activity name. The events within a process instance are ordered. Therefore, a process instance is often represented as a trace over a set of activities. In addition, events can have attributes such as timestamps, associated resources (e.g., the person executing the activity), transactional information (e.g., start, complete, suspend, and so on), and data attributes (e.g., amount or type of customer). For a more formal definition of event logs used in process mining, the reader is referred to [6]. Fig. 2 shows a fragment of an example log. Event logs like in Fig. 2 are completely standard in the process mining community and event log formats such as MXML [7] and XES [8] are used.

The topics in process mining can be broadly classified into three categories: 1) discovery; 2) conformance; and 3) enhancement [6]. Process discovery deals with the discovery of models from event logs. These models may describe control flow, organizational aspects, time aspects, and so on. For example, there are dozens of techniques that automatically construct process models (e.g., Petri nets or BPMN models) from event logs [6]. Fig. 2 shows the basic idea of process discovery. An event log containing detailed information about events is transformed into a multiset of traces

$$L = [abcdjklm, aefjkln, abgchdjkln, \dots]$$

Process discovery techniques are able to discover process models such as the Petri net shown in Fig. 2. Conformance deals with comparing an a priori process model with the observed behavior as recorded in the log and aims at detecting inconsistencies/deviations between a process model and its corresponding execution log. In other words, it checks for any violation between *what was expected to happen* and *what actually has happened*. Enhancement deals with extending or improving an existing model based on information about the process execution in an event log. For example, annotating a process model with performance data to show bottlenecks, throughput times, and so on.

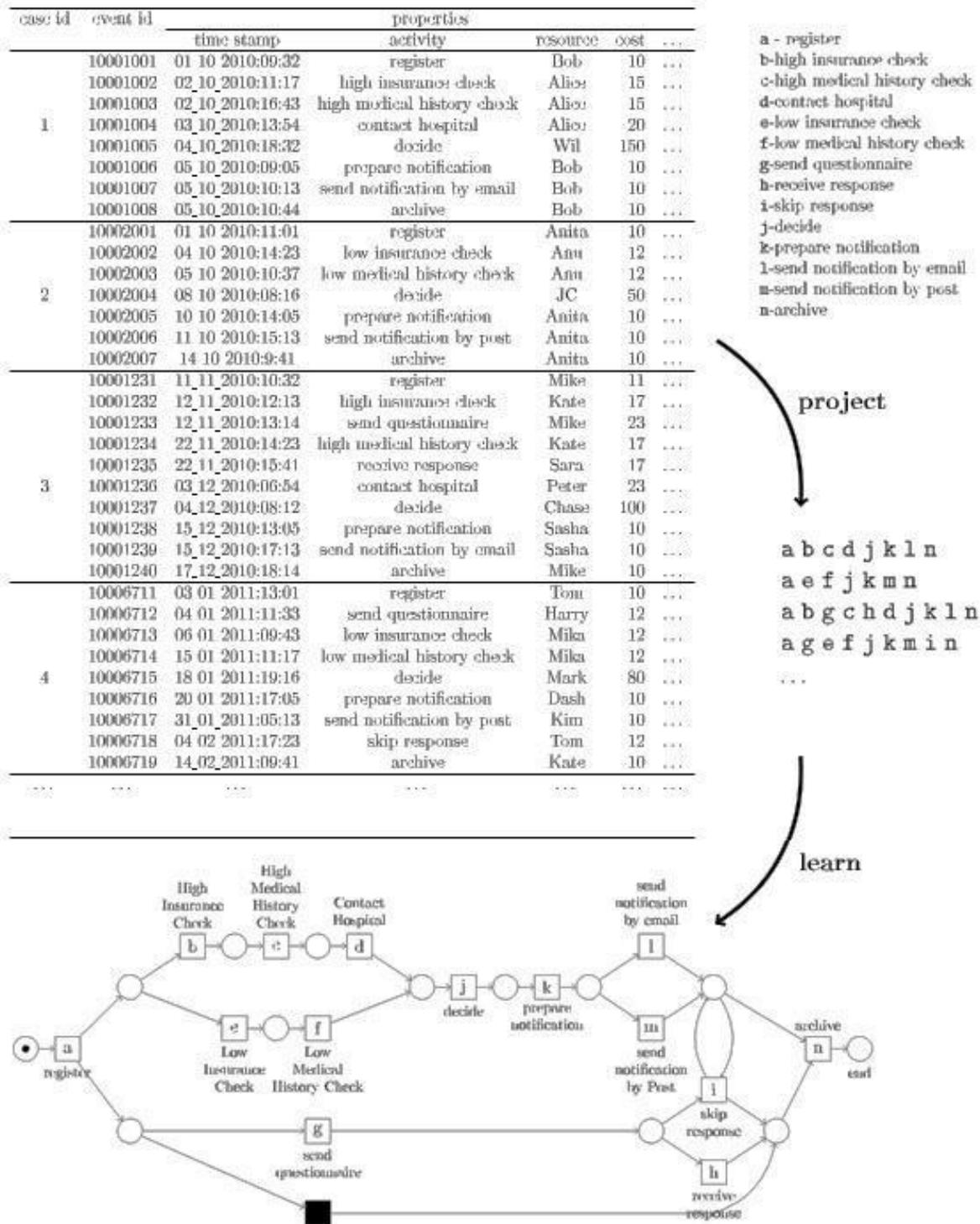


Fig. 2: Process discovery aims to learn a process model (in this case a Petri net) from event logs. An event log consists of events related to cases and referring to activities. To discover control flow, traces are projected onto activity names.

Being a relatively young research discipline, several process mining challenges remain to be addressed. The process mining manifesto [9] lists 11 challenges. The fourth challenge is dealing with concept drift and, thus far, a little work has been done on this highly relevant topic [10], [11].

B. Concept Drift

Concept drift [12] in machine learning and data mining refers to situations when the relation between the input data and the target variable, which the model is trying to predict,

changes over time in unforeseen ways. Therefore, the accuracy of the predictions may degrade over time. To prevent that, predictive models need to be able to adapt online, i.e., to update themselves regularly with new data.

The setting is typically looped over an infinite data stream as follows:

- 1) receive new data;
- 2) make a prediction;
- 3) receive feedback (the true target value); and

4) update the predictive model.

While operating under such circumstances, predictive models are required:

- 1) to react to concept drift (and adapt if needed) as soon as possible;
- 2) to distinguish drifts from once-off noise and adapt to changes, but be robust to noise; and
- 3) to operate in less than data arrival time and use limited memory for storage. In this setting, many adaptive algorithms have been developed (e.g., overviews [13], [14]).

Concept drift is a relatively young research topic that has gained popularity in data mining and machine learning communities in the last 10 years. Concept drift research primarily has been focusing on two directions: 1) how to detect drifts (changes) online (e.g., [15]–[20]) and 2) how to keep predictive models up to date (e.g., [21]–[23]). Concept drift has been shown to be important in many applications (e.g., [24]–[26]). The basis for drift detection could be a raw data stream, a stream of prediction errors, and, more rarely, a stream of predictions or a stream of updated model parameters. Two types of concept drift detection approaches have been used: monitoring evolution of a stream [15], [17] or comparing data distributions in two time windows [16], [18]. The cumulative sum (CUSUM) approach [27] is a representative sequential analysis technique for change detection, different extensions to which have been proposed. One notable example is computational intelligence-based CUSUM or CI-CUSUM [19] that aims to detect a non-stationarity condition by monitoring a multidimensional vector, i.e., multiple features. Adaptive windowing [16] is a representative approach for online change detection using an adaptive size sliding detection window. In this paper, we consider offline change detection and its localization and therefore focus on studying what features to monitor and how to identify when these characteristics change.

III. CHARACTERIZATION OF CHANGES IN BUSINESS PROCESSES

In this section, we discuss the various aspects of process change. Initially, we describe change perspectives (control flow, data, and resource). Then, the different types of drift (sudden, gradual, recurring, periodic, and incremental) are discussed.

A. Perspectives of Change

There are three important perspectives in the context of business processes: 1) control flow; 2) data; and 3) resource. One or more of these perspectives may change over time.

1) Control flow/behavioral perspective:

This class of changes deals with the behavioral and structural changes in a process model. Just like the design patterns in software engineering, there exist change patterns capturing the common control-flow changes [29]. Control-flow changes can be classified into operations such as insertion, deletion, substitution, and reordering of process fragments. For example, an organization which used to collect a fee after processing and acceptance of an application can now change their process to enforce payment of that fee before processing an application. Here, the reordering change pattern had been applied on the

payment and the application processing process fragments. As another example, with the addition of new product offerings, a choice construct is inserted into the product development process of an organization. In the context of PAISs, various control-flow change patterns have been proposed in [28], [29]. Most of these control-flow change patterns are applicable to traditional information/workflow systems as well.

Sometimes, the control-flow structure of a process model can remain intact but the behavioral aspects of a model change. For example, consider an insurance agency that classifies claims as high or low depending on the amount claimed. An insurance claim of e1000 which would have been classified as high last year is categorized as a low insurance claim this year because of the organization's decision to increase the claim limit. The structure of the process remains intact but the routing of cases changes.

2) Dataperspective:

This class of changes refers to the changes in the production and consumption of data and the effect of data on the routing of cases. For example, it may no longer be required to have a particular document when approving a claim.

3) Resource perspective:

This class deals with the changes in resources, their roles, and organizational structure, and their influence on the execution of a process. For example, there could have been a change pertaining to who executes an activity. Roles may change and people may change roles. As another example, certain execution paths in a process could be enabled (disabled) upon the availability (nonavailability) of resources. Furthermore, resources tend to work in a particular manner and such working patterns may change over time, e.g., a resource can have a tendency of executing a set of parallel activities in a specific sequential order. Such working patterns could be more prominent when only few resources are available; the addition of new resources can remove this bias.

B. Nature of Drifts

With the duration for which a change is active, we can classify changes into momentary and permanent. Momentary changes are short lived and affect only a very few cases, whereas permanent changes are persistent and stay for a while [31]. In this paper, we focus on permanent changes as momentary changes often cannot be discovered because of insufficient data. Momentary changes correspond to the notion of outliers/noise in data mining. Changes are perceived to induce a drift in the concept (process behavior). As shown in Fig. 3, we identify four classes of drifts.

1) Suddendrift:

This corresponds to a substitution of an existing process M1 with a new process M2, as shown in Fig. 3(a). M1 ceases to exist from the moment of substitution. In other words, all cases (process instances) from the instant of substitution emanate from M2. This class of drifts is typically seen in scenarios such as emergencies, crisis situations, and change of law. As an example, a new regulation by the finance ministry of India mandates all

banks to procure and report the customer's personal account number in their transactions.

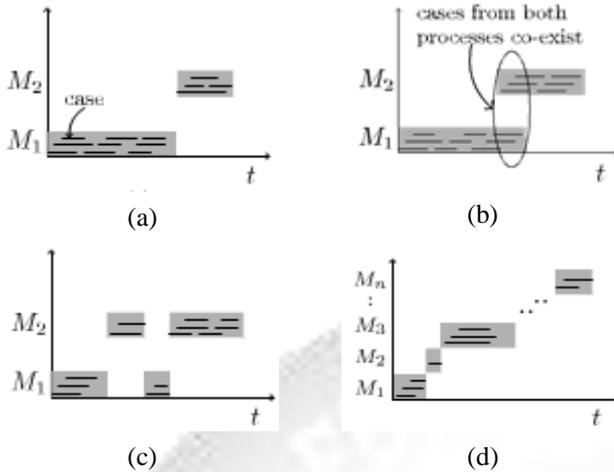


Fig. 3: Different types of drifts. x-axis: time. y-axis: process variants. Shaded rectangles: process instances. (a) Sudden drift. (b) Gradual drift. (c) Recurring drift. (d) Incremental drift.

2) Gradual drift:

This refers to the scenario, as shown in Fig. 3(b) where a current process M_1 is replaced with a new process M_2 . Unlike the sudden drift, here both processes coexist for some time with M_1 discontinued gradually. For example, a supply chain organization might introduce a new delivery process. This process is, however, applicable only for orders taken henceforth. All previous orders still have to follow the former delivery process.

3) Recurring drift:

This corresponds to the scenario where a set of processes reappear after some time (substituted back and forth), as shown in Fig. 3(c). It is quite natural to observe such a phenomenon with processes having a seasonal influence. For example, a travel agency might deploy a different process to attract customers during Christmas period. The recurrence of processes may be periodic or nonperiodic. An example of a nonperiodic recurrence is the deployment of a process subjected to market conditions. The point of deployment and the duration of deployment are both dependent on external factors (here, the market conditions). Periodic drifts may be caused by seasonal effects, e.g., during the summer holidays there tends to be less demand and fewer resources thus influencing the process.

4) Incremental drift:

This refers to the scenario where a substitution of process M_1 with M_N is done via smaller incremental changes, as shown in Fig. 3(d). This class of drifts is more pronounced in organizations adopting an agile BPM methodology and in processes undergoing sequences of quality improvements (most total quality management) initiatives are examples of incremental change [39].

Recurring and incremental drifts in Fig. 3 are shown as discrete sudden changes. These two types of concept drift, however, can also be gradual. Similar categorization of drifts have been proposed in [40] in the context of machine learning. Drifts in [40] are further classified based on the severity of change into severe (and intersected). The categories of severity, as defined in

[40], are too coarse to be applied to business process changes. Nonetheless, the degree of severity in process changes and their impact on dealing with concept drifts is an interesting topic for further research. In the rest, we propose approaches to detect potential control-flow changes in a process manifested as sudden/gradual drifts over a period. Detecting drifts in the other perspectives are beyond the scope of this paper. In addition, as already shown in Fig. 1, we focus on offline concept drift analysis (although our techniques can easily be adapted to the online setting). In practice, a mixture of any or all of the drifts may happen.

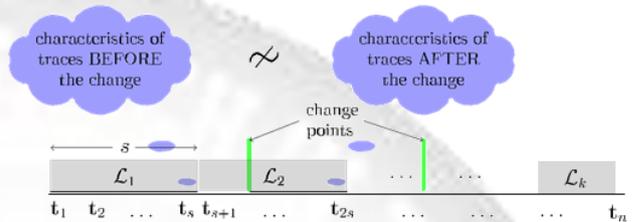


Fig. 4: Event log visualized as a time series of traces along with change points. The basic premise of change (point) detection is that characteristic differences exist in the traces before and after the change

IV. BASIC IDEA OF DRIFT DETECTION IN EVENT LOGS

In this section, we present the basic idea for the detection of changes by analyzing event logs. Initially, we introduce the notations used in this paper.

- 1) A is the set of activities. A^+ is the set of all nonempty finite sequences of activities from A .
- 2) A process instance (i.e., case) is described as a trace over A , i.e., a finite sequence of activities. Examples of traces are $abcd$ and $abbbad$.
- 3) Let $t = t(1)t(2)t(3) \dots t(n) \in A^+$ be a trace over A . $|t| = n$ is the length of the trace t . $t(k)$ is the k^{th} activity in the trace and $t(i, j)$ is the continuous subsequence of t that starts at position i and ends at position j . $t^i = (i, |t|)$ represents the suffix of t that begins at position i .
- 4) An event log, L , corresponds to a multiset (or bag) of traces from A^+ . For example, $L = [abcd, abcd, abbbad]$ is a log consisting of three cases. Two cases follow trace $abcd$ and one case follows trace $abbbad$.
- 5) \mathbb{N} , \mathbb{N}_0 , and \mathbb{R}^+ are the set of all natural numbers, the set of all natural numbers including zero, and the set of all positive real numbers including zero, respectively.

We can consider an event log L as a time series of traces (traces ordered based on the timestamp of the first event). Fig. 4 shows such a perspective on an event log along with change points in the sudden drift scenario.

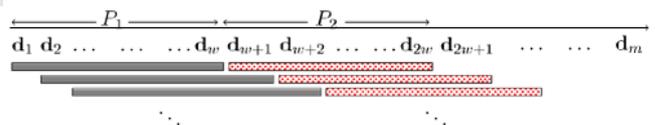


Fig. 7:

Basic idea of detecting drifts using hypothesis tests. The dataset of feature values is considered as a time series for hypothesis tests. P_1 & P_2 are two populations of size w .

V. HYPOTHESIS TESTS FOR DRIFT DETECTION

An event log can be transformed into a data stream/sequence D by choosing one of the feature sets defined in the previous section. The dataset D of feature values can be considered as a series of m values, as shown in Fig. 7. Each $\mathbf{d}_i \in D$ corresponds to the feature value(s) for a trace (or sublog) and can be a scalar or a vector (depending on the choice of feature).⁵ Comparing with Fig. 4, $m = n$ or $m = k$ depending on whether the feature values are recomputed for each trace or for each sublog, respectively. As mentioned earlier, we expect a characteristic difference in the manifestation of feature values in the traces (sublogs) before and after the change points with the difference being more pronounced at the boundaries. To detect this, we can consider a series of successive populations of values (of size w) and investigate if there is a significant difference between two subsequent populations. The premise is that differences are expected to be perceived at change points provided appropriate characteristics of the change are captured as features. A moving window of size w is used to generate the populations. Fig. 7 shows a scenario where two populations $P_1 = \mathbf{d}_1, \mathbf{d}_2, \dots, \mathbf{d}_w$ and $P_2 = \mathbf{d}_{w+1}, \mathbf{d}_{w+2}, \dots, \mathbf{d}_{2w}$ of size w are considered. In the next iteration, the populations correspond to $P_1 = \mathbf{d}_2, \mathbf{d}_3, \dots, \mathbf{d}_{w+1}$ and $P_2 = \mathbf{d}_{w+2}, \mathbf{d}_{w+3}, \dots, \mathbf{d}_{2w+1}$. Given a dataset of m values, the number of population pairs (iterations) will be $m - 2w + 1$.

We propose the use of statistical hypothesis testing to discover these change points. Hypothesis testing is a procedure in which a hypothesis is evaluated on a sample data. One of the important uses of hypothesis testing is to evaluate and compare groups of data. Numerous varieties of hypothesis tests exist [43]. The choice of a particular test is largely dependent on the nature of the data and the objectives of an experiment. For example, hypothesis tests can be classified into parametric and nonparametric tests. Parametric tests assume that the data have a particular distribution, e.g., normal, whereas the nonparametric tests do not make any assumption with regard to the data distribution. Because we do not know the *a priori* distribution of the feature values in an event log, we consider only nonparametric tests. Another perspective of classification is based on the number of samples (populations) on which the hypothesis is defined. We can classify the hypothesis tests into 1) one-sample; 2) two-sample; and 3) multi-sample tests. Because we need to analyze two populations for detecting drifts, we are interested in two-sample hypothesis tests. Another classification of hypothesis tests is concerned with the dimensionality of each data element in a sample. Tests dealing with scalar data elements are called univariate tests while those dealing with vector data elements are called multivariate tests. If only a particular activity or activity pair is considered, then every data item $\mathbf{d} \in D$ is a scalar value corresponding to the trace/sublog i . If we, however, consider a set of activities or activity pairs, then each data item is a vector. Therefore, we need to consider both univariate and multivariate hypothesis tests.

We will use the univariate two-sample Kolmogorov-Smirnov test (KStest) and Mann-Whitney U test (MWtest) as hypothesis tests for univariate data, and the two-sample Hotelling T^2 test for multivariate data. The KStest evaluates

the hypothesis “Do the two independent samples represent two different cumulative frequency distributions?” whereas the MWtest evaluates the hypothesis “Do the two independent samples have different distributions with respect to the rank ordering of the values?”. The multivariate Hotelling T^2 test is a generalization of the t -test and evaluates the hypothesis “Do the two samples have the same mean pattern?”. All of these tests yield a significance probability assessing the validity of the hypothesis on the samples. We refer [43] for a classic introduction to various hypothesis tests.

VI. FRAMEWORK

We propose the framework shown in Fig. 8 for analyzing concept drifts in process mining. The framework identifies the following steps:

1) Feature extraction and selection:

This step pertains to defining the characteristics of the traces in an event log. In this paper, we have defined four features that characterize the control-flow perspective of process instances in an event log. Depending on the focus of analysis, we may define additional features, e.g., if we are interested in analyzing changes in organizational/resource perspective, we may consider features derived from social networks as a means of characterizing the event log. In addition to feature extraction, this step also involves feature selection. Feature selection is important when the number of features extracted is large. We may consider dimensionality reduction techniques [44], [45] such as PCA [46] or random projection [47] to deal with high dimensionality.

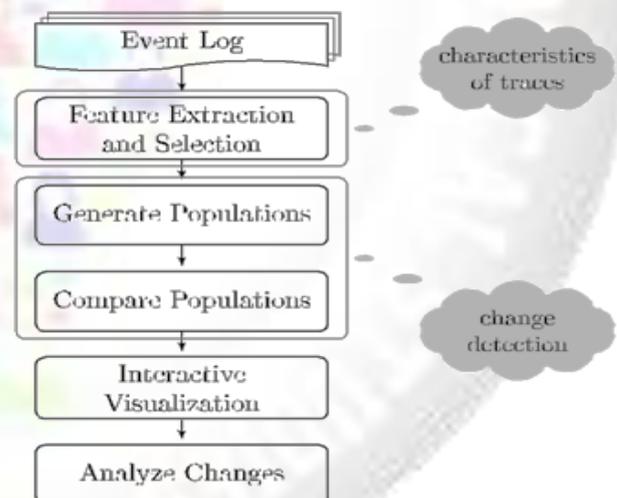


Fig. 8 : Framework for handling concept drifts in process mining.

2) Generate populations:

An event log can be transformed into a data stream based on the features selected in the previous step. This step deals with defining the sample populations for studying the changes in the characteristics of traces. Different criteria/scenarios may be considered for generating these populations from the data stream. In Section VII, we have considered nonoverlapping, continuous, and fixed-size windows for defining the populations. We may also consider, for example, noncontinuous windows (there is a gap between two populations), adaptive windows (windows

can be of different lengths) [16], and soon, which are more appropriate for dealing with gradual and recurring drifts.

3) *Compare populations:*

Once the sample populations are generated, the next step is to analyze these populations for any change in characteristics. In this paper, we advocate the use of statistical hypothesis tests for comparing populations. The null hypothesis in statistical tests states that distributions (or means, or standard deviations) of the two sample populations are equal. Depending on desired assumptions and the focus of analysis, different statistical tests can be used.

4) *Interactive visualization:*

The results of comparative studies on the population of trace characteristics can be intuitively presented to an analyst. For example, the significance probabilities of the hypothesis tests can be visualized as a drift plot. Troughs in such a drift plot signify a change in the significance probability, thereby implying a change in the characteristics of traces.

5) *Analyze changes:*

Visualization techniques such as the drift plot can assist in identifying the change points. Having identified that a change had taken place, this step deals with techniques that assist an analyst in characterizing and localizing the change and in discovering the change process. The framework can be used for designing new change detection approaches.

VII. IMPLEMENTATION

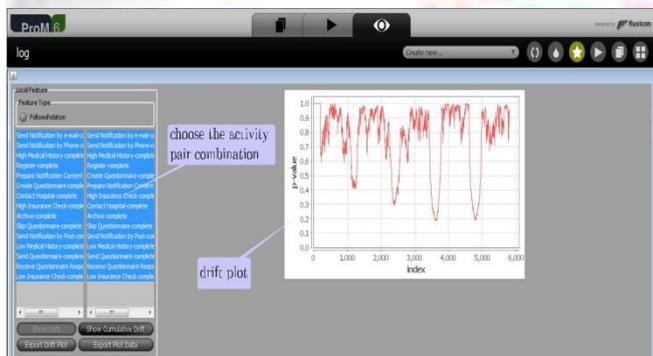


Fig. 9: Visualization of the drift plot in the concept drift plug-in in ProM.

The concepts presented in this paper have been realized as the concept drift plug-in in the ProM⁶ framework. ProM is a plug-able environment for process mining envisioned to provide a common basis for all kinds of process mining techniques ranging from importing, exporting, and filtering event logs (process models) to analysis and visualization of results. Over years, ProM has emerged to be the de facto standard for process mining. The concept drift plug-in implements all of the steps in the proposed framework and can be easily extended with additional elements (e.g., new features can be easily added). The plug-in supports visualization of the significance probability for the hypothesis tests as a drift plot. Fig. 9 shows a drift plot from the plug-in.

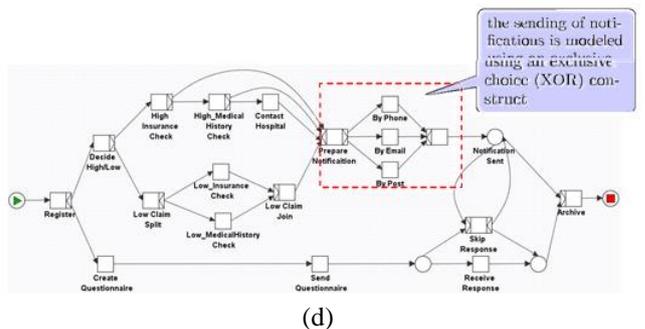
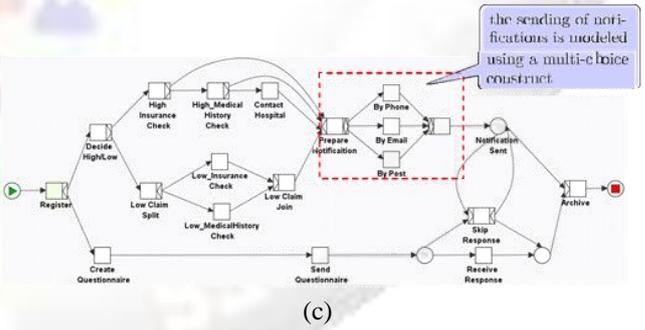
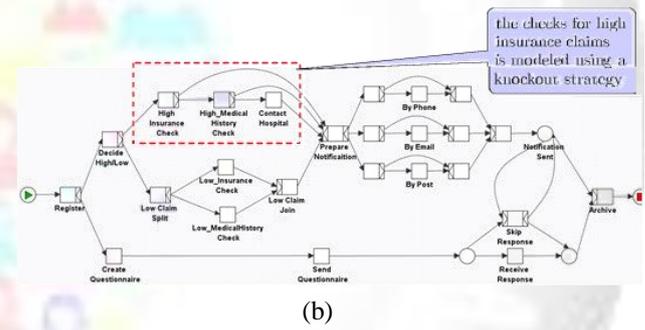
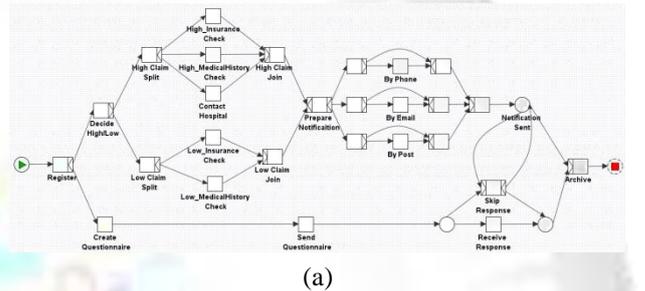
VIII. EXPERIMENTAL RESULTS AND DISCUSSION

Now, we put the ideas proposed for handling concept drifts in practice. Initially, we will illustrate the effectiveness of the

proposed approaches using a synthetic example of an insurance claim process and later discuss the results from a real-life case study in a large Dutch municipality.

A. *Synthetic Log-Insurance Claim Process*

This process corresponds to the handling of health insurance claims in a travel agency. Upon registration of a claim, a general questionnaire is sent to the claimant. In parallel, a registered claim is classified as high or low. For low claims, two independent tasks: 1) check insurance and 2) check medical history need to be executed. For high claims, three tasks need to be executed: 1) check insurance; 2) check medical history; and 3) contact doctor/hospital for verification. If one of the checks show that the claim is not valid, then the claim is rejected; otherwise, it is accepted. A cheque and accepted decision letter is prepared in cases where a claim is accepted while a rejection decision letter is created for rejected claims. In both cases, a notification is sent to the claimant.



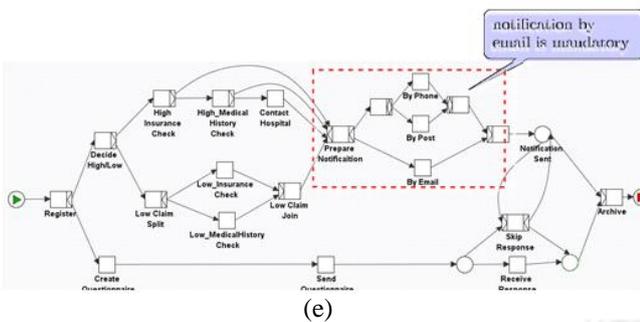


Fig. 10 : Variants of an insurance claim process of a travel agency represented in YAWL notation. Dashed rectangles: regions of change from its previous model. (a) Model 1. (b) Model 2. (c) Model 3. (d) Model 4. (e) Model 5.

Three modes of notification are supported by: 1) email; 2) telephone (fax); and 3) postal mail. The cases should be archived upon notifying the claimant. This can be done with or without the response for the questionnaire. The decision of ignoring the questionnaire, however, can only be made after a notification is sent. The case is closed upon completion of archiving task. Fig. 10 shows five variants of this process represented in YAWL [48] notation. Dashed rectangles: a change has been done in the process model with respect to its previous variant. The changes can have various reasons. For example, in Fig. 10(a), the different checks for high insurance claims are modeled using a parallel (AND) construct. A claim, however, can be rejected if any one of the checks fail. In such cases, the time and resources spent on other checks go waste. To optimize this process, the agency can decide to enforce an order on these checks and proceed on checks only if the previous check results are positive. In other words, the process is modified with a *knockout* strategy [49] adopted for the process fragment involving the different checks for high insurance claims, as shown in Fig. 10(b). As another example, the OR-construct pertaining to the sending of notification to claimants in Fig. 10(c) has been modified to an exclusive-or (XOR) construct in Fig. 10(d). The organization could have taken a decision to reduce their workforce as a cost-cutting measure. Because of the availability of limited resources, they would like to minimize the redundancy of sending the notification through different modes of communication and restrict it to only one of the modes. Considering an event log containing cases that belong to such a mix of process variants, the objective of change point detection is to detect when the processes have changed. In this section, we will illustrate the handling of concept drifts in the context of sudden and gradual drifts. We have modeled each of these five process variants in CPN tools [50] and simulated 1200 traces for each model.

1) Sudden Drift Change (Point) Detection:

To simulate the sudden drift phenomenon, we created an event log L consisting of 6000 traces by juxtaposing each set of the 1200 traces. The event log contains 15 activities or event classes (i.e., $|A| = 15$) and 58783 events (which is the total number of events in the log for all the traces). Given this event log L , our first objective is to detect the four change points pertaining to these five process variants, as shown in Fig. 11(a). Global features can be applied only at the log level; to facilitate this, we have split the log into 120 sublogs using a split size of 50 traces.

In this scenario, the four change points corresponding to the five process variants are, as shown in Fig. 11(b). We have computed the follows RCo of all 15 activities thereby generating a multivariate vector of 45 features for each sublog. We have applied the Hotelling T^2 hypothesis test on this multivariate dataset using a moving window population of size, $w=10$. For this hypothesis test, we have randomly chosen 12 of the 45 features with a 10-fold cross validation. Fig. 12(a) shows the average significance probability of the Hotelling T^2 test for the 10 folds on this feature set. The troughs in the plot signify that there is a change in the distribution of the feature values in the log. In other words, they show that there is a drift (change) in the concept, which here corresponds to the process. It is interesting to observe that the troughs are observed around indexes 24, 72, and 96 which are indeed the points of change (remember that, we have split the log into 120 sublogs with the change points at indexes 24, 48, 72, and 96). The change at index 48 corresponding to the transition from M_2 to M_3 could not be uncovered using this feature set because the RCs would be alike for logs generated from these two process variants.

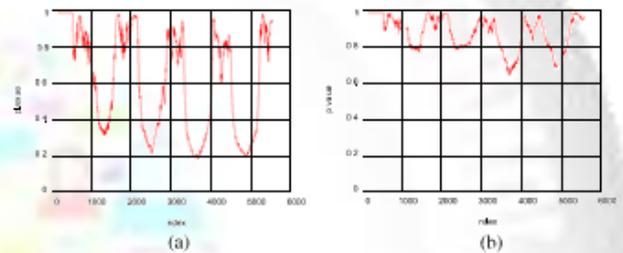


Fig. 13:

Average significance probability (overall activity pairs) of KStest on the J measure and WC feature set estimated for each trace. x -axis: trace index. y -axis: significance probability of the test. Troughs: change points. Vertical grid lines: actual (sudden) change points. (a) J -measure. (b) WC.

corresponding to each activity pair in $A \times A$. Fig. 13(a) shows the average significance probability of KStest on all activity pairs, whereas Fig. 13(b) shows the average significance probability of KStest on all activity pairs using the WC feature set. We can observe that significant troughs are formed at indexes 1200, 2400, 3600, and 4800. These are indeed the points where the model has been changed.

• Influence of Population Size:

It is imperative to note that the goodness of the result of hypothesis tests depends on the population size. The statistical analysis assumes that each population is independent. A good population size is largely dependent on the application and the focus of analysis. To study the influence of population size, we have considered the J measure for every pair of activities and the univariate KS test for change point detection. Fig. 14 shows the results for varying sizes of the population. We observe a lot of noise for small populations and the drift tends to be smooth as the population size increases. This can be attributed to the fact that as the population size increases (i.e., as we consider more cases), the variability in the nature of cases reduces and attains a stability, e.g., there can be a flux of low-insurance claims initially and after a certain time the proportion stabilizes.

2) *SuddenDrift ChangeLocalization:*

Oursecondobjective inhandling concept driftsisthatofchangelocalization. To localizethechanges (identifytheregions ofchange),weneed to consider activity pairs individually or subsets of activitypairs. For example, the change from M1 to M2 is localized in the region pertaining to high insurance claim checks. We expect characteristic changes in features pertaining to these activities and other activities related to these activities. For example, in M1, the activities High Medical History Check and Contact Hospital always follow the activity Register whenever a claim is classified as high. In contrast, in M2, these activities need not always follow Register because both these activities are skipped if High Insurance Check fails while Contact Hospital is skipped if High Medical History Check fails. During simulation, we have set the probability of success of a check to 90%. We have considered the WC feature for the activity relation Contact Hospital follows Register on a window length of $l = 10$ in each trace separately. Fig. 15(a) shows the significance probability of the univariate KS test using a population size of $w = 400$ on this feature. We can observe that one dominant trough is formed at index 1200 showing that there exists a change in the region between Register and Contact Hospital. No subsequent changes with respect to this activity pair can be noticed, which is indeed the case in the sequence of models used.

the accuracy of the proposed framework in handling gradual drifts. Recall that in gradual drifts, one concept fades gradually while the other takes over. This phenomenon of gradual change can be modeled in many ways. In this paper, we consider the scenario where the change is linear between two sources, as shown in Fig. 16(a). In this figure, we observe the fading of one concept M1 and the taking over of another concept M2 happen linearly. Within this setup, we can alter the extent to which the two concepts coexist. For the insurance claim example, we generated two event logs exhibiting gradual drifts by varying the duration of change. In the first case, the process variants M1 and M2 coexist between trace indexes 1000 and 1400, the variants M2 and M3 coexist between indexes 2200 and 2600, the variants M3 and M4 coexist between indexes 3400 and 3800, and the variants M4 and M5 coexist between indexes 4600 and 5000, as shown in Fig. 16(b). The point of cross over is still retained at indexes 1200, 2400, 3600, and 4800.

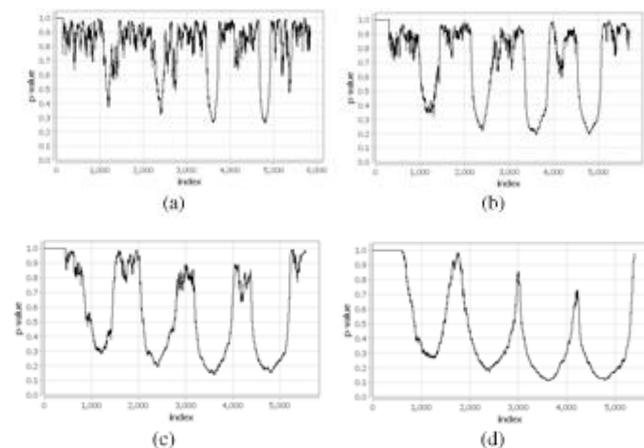


Fig. 14 : Average significance probability (over all activity pairs) for different population sizes of KS test on the J measure estimated over all activity pairs in each trace. x-axis: trace index. y-axis: significance probability of the test. Troughs: change points. Vertical grid lines: actual (sudden) change points. (a) $w = 150$, (b) $w = 300$, (c) $w = 450$, and (d) $w = 600$.

B. *Real-LifeLog:AdvertisementPermitProcess ofaDutchMunicipality*

Thesyntheticeventdatacreatedthroughsimulationallow ustocomparethecontrolled (ground truth)changeswiththe detectedchanges.We,however,havealsoapplied ourconcept driftanalysisstechniques tovariousreal-lifeventlogs.Here, wereportonacasestudywhereweanalyzedconceptdriftsinprocesseswithinalargeDutchmunicipality.Municipalities areinterestedinobtaininginsightsintotheirprocesses,e.g., thewaytheyareplannedtoexecutedvis-a-vistheway

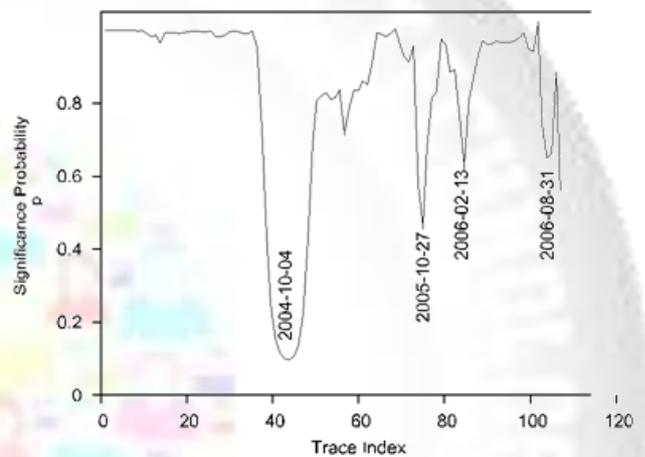


Fig. 19. Average significance probability (over all activity pairs) of KS test on J measure. The population size for the KS test is 10. There are four troughs signifying a change in behavior.

follows.

- 1) Upon submission of an application, the municipality acknowledges the receipt of documents and (optionally) tests for its completeness.
- 2) Then, the municipality proceeds with a follow-up procedure that verifies whether the application and submitted documents are in compliance with the regulations.
- 3) With the investigations, then the municipality makes a decision on the application and informs the applicant with the decision along with a fee letter.
- 4) Finally, the municipality registers the advertisements placed and enforces them.

Fig. 20(b) shows the process model discovered using the Heuristic miner [56] on the event log L . The figure highlights regions that differ from the process model in Fig. 20(a). There are two changes in this model with respect to the previous one. The first change is related to the checking for completeness of the registered documents. In the initial process model [Fig. 20(a)], this check was not mandatory (only two of the 43 applications were checked for completeness). The municipality changed this process by making the checks mandatory before proceeding. The second change is the introduction of a new

activity 'End procedure; enforcement is next,' as shown in Fig. 20(b). The initial process model had only the activity 'End procedure, possibly choose enforcement' whereas the new model has both these activities. Similar changes have been observed in the rest of the models. We do not provide them here because of space constraints. For a more detailed discussion on this case study refer [57].

C. Time Complexity

In this section, we assess the time complexity of the proposed approach. The feature extraction is dependent on the size of the event log (i.e., the number of events). The feature values for all of the features proposed in this paper can be extracted in linear time with respect to the size of the event log. Fig. 21(a) shows the average time along with 95% confidence intervals (over five independent runs) for extracting the J measure

IX. CONCLUSION

In this paper, we have introduced the topic of concept drift in process mining, i.e., analyzing process changes based on event logs. We proposed feature sets and techniques to effectively detect the changes in event logs and identify the regions of change in a process. Our initial results show that heterogeneity of cases arising because of process changes can be effectively dealt with by detecting concept drifts. Once change points are identified, the event log can be partitioned and analyzed. This is the first step in the direction of dealing with changes in any process monitoring and analysis efforts. We have considered changes only with respect to the control-flow perspective manifested as sudden and gradual drifts. Therefore, our analysis should only be observed as the starting point for a new subfield in the process mining domain and there are lots of challenges that still need to be addressed. Some of these challenges include.

- 1) Change-pattern specific features:
- 2) Feature selection:
- 3) Recurring drifts:
- 4) Holistic approaches:
- 5) Change process discovery:
- 6) Sample complexity:
- 7) Online (on-the-fly) drift detection:

REFERENCES

- [1] (2010). All-in-one Permit for Physical Aspects: (Omgevingsvergunning) in a Nutshell [Online]. Available: <http://www.answersforbusiness.nl/regulation/all-in-one-permit-physical-aspects>
- [2] United States Code. (2002, Jul.). Sarbanes-Oxley Act of 2002, PL 107-204, 116 Stat 745 [Online]. Available: <http://files.findlaw.com/news.findlaw.com/cnn/docs/gwbush/sarbanesoxley072302.pdf>
- [3] W. M. P. van der Aalst, M. Rosemann, and M. Dumas, "Deadline-based escalation in process-aware information systems," *Decision Support*
- [4] SyMst., Dvuoml. a4s3, Wn.oM. 2. Pp. pv.an49d2e-r5A1a11st2, 0a1n1d. A. H. M. Ter Hofstede, *Process-Aware Information Systems: Bridging People and Software Through Process Technology*. New York, NY, USA: Wiley, 2005.
- [5] W. M. P. van der Aalst and K. M. van Hee, *Workflow Management: Models, Methods, and Systems*. Cambridge, MA, USA: MIT Press, 2004.
- [6] W. M. P. van der Aalst, *Process Mining: Discovery, Conformance and Enhancement of Business Processes*. New York, NY, USA: Springer-Verlag, 2011.
- [7] B. F. van Dongen and W. M. P. van der Aalst, "A meta model for process mining data," in *Proc. CAiSE Workshops (EMOI-INTEROP Workshop)*, vol. 2. 2005, pp. 309–320.
- [8] C. W. Günther, (2009). XES Standard Definition [Online]. Available: <http://www.xes-standard.org>
- [9] F. Daniel, S. Dustdar, and K. Barkaoui, "Process mining manifesto," in *BPM 2011 Workshops*, vol. 99. New York, NY, USA: Springer-Verlag, 2011, pp. 169–194.
- [10] R. P. J. C. Bose, W. M. P. van der Aalst, I. Žliobaite, and M. Pechenizkiy, "Handling concept drift in process mining," in *Proc. Int. CAiSE*, 2011 pp. 391–405
- [11] J. Carmona and R. Gavalda, "Online techniques for dealing with concept drift in process mining," in *Proc. Int. Conf. IDA*, 2012, pp. 90–102.
- [12] J. Schlimmer and R. Granger, "Beyond incremental processing: Tracking concept drift," in *Proc. 15th Nat. Conf. Artif. Intell.*, vol. 1. 1986, pp. 502–507.
- [13] A. Bifet and R. Kirkby. (2011). *Data Stream Mining: A Practical Approach*, University of Waikato, Waikato, New Zealand [Online]. Available: <http://www.cs.waikato.ac.nz/~abifet/MOA/StreamMining.pdf>
- [14] I. Žliobaite, "Learning under concept drift: An Overview," *CoRR*, vol. abs/1010.4784, 2010 [Online]. Available: <http://arxiv.org/abs/1010.4784>
- [15] J. Gama, P. Medas, G. Castillo, and P. Rodrigues, "Learning with drift detection," in *Proc. SBIA*, 2004, pp. 286–295.
- [16] A. Bifet and R. Gavalda, "Learning from time-changing data with adaptive windowing," in *Proc. 7th SIAM Int. Conf. Data Mining (SDM)*, 2007, pp. 443–448.
- [17] G. J. Ross, N. M. Adams, D. K. Tasoulis, and D. J. Hand, "Exponentially weighted moving average charts for detecting concept drift," *Pattern Recognit. Lett.*, vol. 33, no. 2, pp. 191–198, 2012.
- [18] K. Nishida and K. Yamauchi, "Detecting concept drift using statistical testing," in *Proc. 10th Int. Conf. Discovery Sci.*, 2007, pp. 264–269.
- [19] C. Alippi and M. Roveri, "Just-in-time adaptive classifiers—Part I: Detecting nonstationary changes," *IEEE Trans. Neural Netw.*, vol. 19, no. 7, pp. 1145–1153, Jul. 2008.
- [20] C. Alippi, G. Boracchi, and M. Roveri, "Just-in-time classifiers for recurrent concepts," *IEEE Trans. Neural Netw. Learn. Syst.*, vol. 24, no. 4, pp. 620–634, Apr. 2013.
- [21] J. Z. Kolter and M. A. Maloof, "Dynamic weighted majority: An ensemble method for drifting concepts," *J. Mach. Learn. Res.*, vol. 8, pp. 2755–2790, Jan. 2007.
- [22] R. Elwell and R. Polikar, "Incremental learning of concept drift in nonstationary environments," *IEEE Trans. Neural Netw.*, vol. 22, no. 10, pp. 1517–1531, Oct. 2011.

- [23] G. Widmer and M. Kubat, "Learning in the presence of concept drift and hidden contexts," *Mach. Learn.*, vol. 23, no. 1, pp. 69–101, Apr. 1996.
- [24] S. J. Delany, P. Cunningham, A. Tsymbal, and L. Coyle, "A case-based technique for tracking concept drift in spam filtering," *Knowl. Based Syst.*, vol. 18, nos. 4–5, pp. 187–195, Aug. 2005.
- [25] A. Tsymbal, M. Pechenizkiy, P. Cunningham, and S. Puuronen, "Handling local concept drift with dynamic integration of classifiers: Domain of antibiotic resistance in nosocomial infections," in *Proc. 19th IEEE Int. Symp. CBMS*, Nov. 2006, pp. 679–684.
- [26] M. Pechenizkiy, J. Bakker, I. Žliobaite, A. Ivannikov, and T. Kärkkäinen, "Online mass flow prediction in CFB boilers with explicit detection of sudden concept drift," *SIGKDD Explorations*, vol. 11, no. 2, pp. 109–116, 2009.
- [27] E. S. Page, "Continuous inspection schemes," *Biometrika*, vol. 41, nos. 1–2, pp. 100–115, 1954.
- [28] N. Mulyar, "Patterns for process-aware information systems: An approach based on colored Petri nets," Ph.D. dissertation, Dept. Comput. Sci., Univ. Technol., Eindhoven, The Netherlands, 2009.
- [29] B. Weber, S. Rinderle, and M. Reichert, "Change patterns and change support features in process-aware information systems," in *Proc. 19th Int.*, 2007, pp. 574–588.
- [30] G. Regev, P. Soffer, and R. Schmidt, "Taxonomy of flexibility in business processes," in *Proc. 7th Workshop BPMDS*, 2006, pp. 1–4.
- [31] H. Schonenberg, R. Mans, N. Russell, N. Mulyar, and W. M. P. van der Aalst, "Process flexibility: A survey of contemporary approaches," in *Proc. Adv. Enterprise Eng. I*, 2008, pp. 16–30.
- [32] K. Ploesser, J. C. Recker, and M. Rosemann, "Towards a classification and lifecycle of business process change," in *Proc. BPMDS*, vol. 8, 2008, pp. 1–9.
- [33] C. W. Günther, S. Rinderle-Ma, M. Reichert, and W. M. P. van der Aalst, "Using process mining to learn from process changes in evolutionary systems," *Int. J. Business Process Integr. Manag.*, vol. 3, no. 1, pp. 61–78, 2008.
- [34] M. van Leeuwen and A. Siebes, "StreamKrimp: Detecting change in data streams," in *Proc. Mach. Learn. Knowl. Discovery Databases*, 2008, pp. 672–687.
- [35] I. Žliobaite and M. Pechenizkiy. (2010). Handling Concept Drift in Information Systems [Online]. Available: http://sites.google.com/site/zliobaite/CD_applications_2010.pdf
- [36] H. Wang, W. Fan, P. S. Yu, and J. Han, "Mining concept-drifting data streams using ensemble classifiers," in *Proc. 9th ACM SIGKDD Int. Conf. Knowl. Discovery Data Mining*, 2003, pp. 226–235.
- [37] D. Brzezinski and J. Stefanowski, "Reacting to different types of concept drift: The accuracy updated ensemble algorithm," *IEEE Trans. Neural Netw. Learn. Syst.*, Apr. 2013, doi: 10.1109/TNNLS.2013.2251352.
- [38] W. M. P. van der Aalst, A. Adriansyah, and B. Dongen, "Replaying history on process models for conformance checking and performance analysis," *WIREs Data Mining Knowl. Discovery*, vol. 2, no. 2, pp. 182–192, 2012.
- [39] M. Hammer, *Beyond Reengineering: How the Process-Centered Organization is Changing Our Work and Our Lives*. New York, NY, USA: Harper business, 1996.
- [40] L. L. Minku, A. P. White, and X. Yao, "The impact of diversity on online ensemble learning in the presence of concept drift," *IEEE Trans. Knowl. Data Eng.*, vol. 22, no. 5, pp. 730–742, May 2010.
- [41] P. Smyth and R. M. Goodman, *Rule Induction Using Information Theory*. Washington, DC, USA: AAAS Press, 1991, pp. 159–176.
- [42] N. M. Blachman, "The amount of information that y gives about X," *IEEE Trans. Inf. Theory*, vol. 14, no. 1, pp. 27–31, Jan. 1968.
- [43] D. Sheskin, *Handbook of Parametric and Nonparametric Statistical Procedures*. London, U.K.: Chapman & Hall/CRC, 2004.
- [44] I. K. Fodor, "A survey of dimensionality reduction techniques," in *Proc. Center Appl. Sci. Comput., Lawrence Livermore Nat. Lab.*, 2002, pp. 1–24.
- [45] I. Guyon and A. Elisseeff, "An introduction to variable and feature selection," *J. Mach. Learn. Res.*, vol. 3, pp. 1157–1182, Mar. 2003.