

Predictive Process Monitoring in Apromore

Ilya Verenich^{1,2}, Stanislav Mõškovski², Simon Raboczi³, Marlon Dumas²,
Marcello La Rosa³, and Fabrizio Maria Maggi²

¹ Queensland University of Technology, Australia

`ilya.verenich@qut.edu.au`

² University of Tartu, Estonia

`{stanislav.myshkovski, marlon.dumas, f.m.maggi}@ut.ee`

³ University of Melbourne, Australia

`{simon.raboczi, marcello.larosa}@unimelb.edu.au`

Abstract. This paper discusses the integration of Nirdizati, a tool for predictive process monitoring, into the Web-based process analytics platform Apromore. Through this integration, Apromore’s users can use event logs stored in the Apromore repository to train a range of predictive models, and later use the trained models to predict various performance indicators of running process cases from a live event stream. For example, one can predict the remaining time or the next events until case completion, the case outcome, or the violation of compliance rules or internal policies. The predictions can be presented graphically via a dashboard that offers multiple visualization options, including a range of summary statistics about ongoing and past process cases. They can also be exported into CSV for periodic reporting or to be visualized in third-parties business intelligence tools. Based on these predictions, operations managers may identify potential issues early on, and take remedial actions in a timely fashion, e.g. reallocating resources from one case onto another to avoid that the case runs overtime.

Keywords: Process mining, predictive monitoring, business process, machine learning

1 Introduction

Predictive Process Monitoring is an emerging paradigm based on the continuous generation of predictions about the future values of user-specified performance indicators of a currently running process execution [4]. In this paradigm, a user defines the type of predictions they are interested in and provides a set of historical execution traces. Based on the analysis of these traces, the idea of predictive monitoring is to continuously provide the user with predictions and estimated values of the performance indicators. Such predictions generally depend both on: (i) the sequence of activities executed in a given case; and (ii) the values of data attributes after each activity execution in the case.

As an example, consider a business process for managing loan applications; the applicant can be advised about the combinations of loan amount and duration that are the most likely to lead to an acceptance of the application, given the history of the application until now and the personal data of the applicant (e.g., age, salary, etc.).

A range of approaches have been proposed in the literature to tackle common predictive process monitoring tasks [2]. However, these approaches have largely remained in the academic domain and have not been widely applied in real-time scenarios where users require a continuous predictive support.

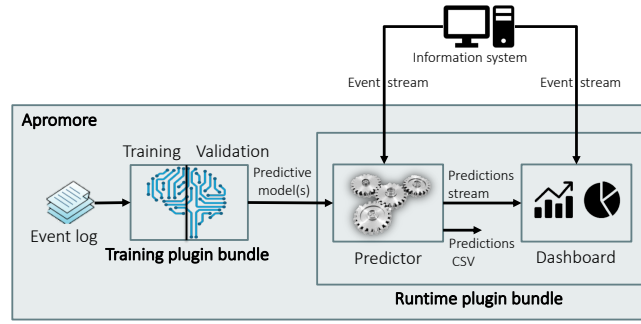


Fig. 1. High-level architecture of the predictive monitoring functionality of Apromore.

To fill this gap, last year we started the *Nirdizati* open-source project, aimed at developing a Web-based tool for the predictive monitoring of business processes. In this paper, we document the developments of *Nirdizati* since its creation, and its integration into the process analytics platform Apromore.

The two core components of *Nirdizati*, namely *Training* and *Runtime*, have been integrated as two bundles (i.e. sets) of plugins into Apromore (Fig. 1). The Training plugin bundle takes as input a business process event log stored in the Apromore repository, and produces one or more predictive models, which can then be deployed to the runtime predictive monitoring environment. Once a model is deployed, the Runtime plugin bundle listens to a stream of events coming from an information system supporting the process, or produced by replaying an event log stored in the repository, and creates a stream of predictions. These predictions can then be visualized in a Web dashboard or exported into a CSV (comma-separated value) file for periodic reporting and to be used within third-party business intelligence tools.

2 Apromore Platform

Apromore is a Web-based advanced process analytics platform, developed by the BPM community under an open-source initiative.¹ Apromore was originally conceived as an advanced process model repository. However, today it offers a wide range of features which go beyond those for managing large process model collections, and include a variety of state-of-the-art process mining techniques. These are techniques for the automated discovery of BPMN models, for the conformance checking of BPMN models against event logs, the replaying of event logs on top of BPMN models, the detection and characterization of process drifts from event logs, the visual analysis of process performance, and many others.

All these features are exposed through a Web portal, and organized according to the phases of the BPM lifecycle: discovery, analysis, redesign, implementation and monitoring. These features can also be accessed as external Web services by third-party BPM software environments, such as ProM (for process mining) and WoPeD (for process modeling and verification).

From a technology viewpoint, Apromore relies on four core technologies: Spring, ZK, OSGi and Eclipse Virgo. Spring provides a simplified management of Java-based

¹ <http://apromore.org>

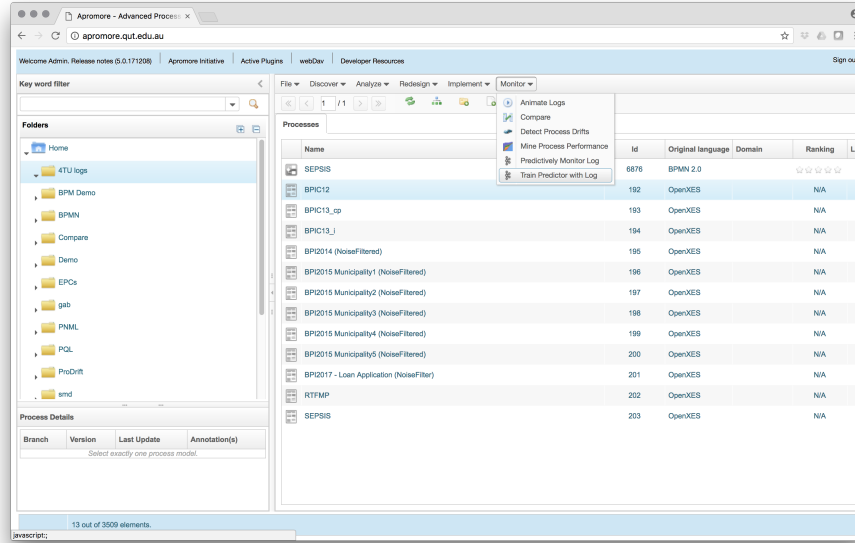


Fig. 2. Apromore's Portal with predictive monitoring functionality highlighted.

enterprise applications through the use of Java annotations and XML configurations. ZK is an AJAX framework used for Apromore's main Web interface (the *Portal*). OSGi provides a flexible framework for managing component dependencies through plugin bundles. Finally, Eclipse Virgo is a Web server based on the OSGi component model.

To equip Apromore with predictive process monitoring capabilities, we have wrapped the two core components of Nirdizati into two OSGi plugin bundles for Apromore: Training and Runtime. Each bundle is a set of OSGi plugins which encapsulate the logic or the user interface (UI) of the various functions offered by Nirdizati. For example, the runtime predictor is a logic plugin, while the runtime dashboard is a portal plugin (UI). These two bundles are accessible from the Monitoring menu of the Apromore Portal (see Figure 2). One can select an event log stored in the repository, and use it to train, tune and test a variety of predictive models, by launching the training plugin bundle. Next, the runtime bundle can be used to stream an event log from the repository, or hook into a live external stream, in order to generate predictions as process cases unfold.

In the next sections, we introduce a working example and use this to describe the functionality of the training and runtime plugins in detail.

3 Running Example

As a running example, throughout this paper we will consider a purchase-to-pay process of an IT vendor. The process starts with lodging a purchase order and ends when requested goods have been supplied.

For use with the Training plugin bundle, we extracted an event log of completed purchased orders, while ongoing orders were fed into the Runtime plugin bundle to make predictions for them. The corresponding event log contains a number of case attributes, such as the total cost of the order, supplier identifier, supplier location, delivery

The screenshot shows a web-based training configuration interface. On the left, there are dropdowns for 'Currently selected log' (set to 'BPI2012W.csv') and 'What do you want to predict?' (set to 'Remaining time'). Below these is a 'Threshold' section with radio buttons for 'Average' and 'Custom' (set to 0.1), and an 'Advanced mode' toggle. A link 'Generate new dataset parameters' is also present. The main area contains three sections: 'Encoding' with checkboxes for 'Last state', 'Frequency', 'Combined', and 'Index based'; 'Bucketing method' with checkboxes for 'Zero', 'State based', 'Clustering', and 'Prefix length based'; and 'Prediction method' with checkboxes for 'Random forest', 'Gradient boosting', 'Decision tree', and 'XGBoost'. Below these is a 'GRADIENT BOOSTING' section with input fields for 'Number of estimators' (300), 'Max features' (0.5), and 'Learning rate' (0.1). At the bottom is a blue 'TRAIN MODEL' button.

Fig. 3. Training configuration screen.

type and the target supply date. Moreover, each event contains identifier of an employee who performed it.

In this process, stakeholders are interested in predicting four variables:

- *Late supply*. It is a boolean variable indicating whether or not the case will be closed before the supply date of that case.
- *Delay Rank* indicating the degree of potential delay, namely “Just in case”, “Mild”, “Moderate”, “Severe”.
- *Next activity* indicating which activity will be performed right after the current one.
- *Remaining time* until case completion.

In order to make more accurate predictions, we performed some basic feature engineering before feeding the log into Apromore. For example, to take into account resource contention, we added the number of currently open cases as an event attribute. Furthermore, instead of the absolute target supply date, we used the number of days relative to the case start date. Intuitively, if this number is lower, there exists a higher probability of missing the deadline (Late Supply).

4 Training Plugin Bundle

The Training plugin bundle provides several algorithms for generating predictive models suitable for different types of predictions. Specifically, it is able to build models for predicting remaining time, the next activity to be performed, whether a case will exceed a specified duration threshold, as well as various static case attributes, for example, the total cost of the order. To this aim, the Training bundle involves two phases: a training and a validation phase. In the former, one or more predictive models are fitted; in the latter, their suitability to the specific dataset is evaluated, so as to support the user in selecting the predictive model that ensures the best results.

The Training bundle is composed of a front-end application (Fig. 3), which allows users to select the prediction methods and to assess the goodness-of-fit of the built models, and a back-end application for the actual training and validation. From the data flow perspective, the back-end application performs several tasks shown in Fig. 4.

Firstly, when a user uploads their log, the tool extracts and categorizes data attributes of the log. In order to properly construct feature vectors from business process traces, the attributes need to be categorized into static case attributes and dynamic event attributes. To this end, for each attribute in the log, we check whether it changes through-

out the case lifetime. On the other hand, each attribute needs to be designated as either numeric or categorical.

In order to that, we count the number of unique levels in the values of each attribute and if it less than a certain threshold, we decide the attribute to be categorical, otherwise it is numeric. These procedures are performed automatically upon the log uploading. Nevertheless, the user is given an option to override the automatic attribute definitions. The resulting definitions are saved in a configuration file in a JSON format (Fig. 5). Secondly, the log is internally split into training and validation set in a 80-20 proportion. The former is used to train the model, while the latter is used to evaluate the predictive power of the model. Next, all traces of a business process need to be represented as fixed-size feature vectors in order to train a predictive model. To this end, several encoding techniques were proposed in [1] and further refined in [3], out of which we support four, namely last state encoding, frequency (aggregation) encoding, combined encoding and lossless index-based encoding.

While some of existing predictive process monitoring approaches train a single classifier on the whole event log, others employ a multi-classifier approach by dividing the prefix traces in the historical log into several buckets and fitting a separate classifier for each such bucket. At run-time, the most suitable bucket for the ongoing case is determined and the respective classifier is applied to make a prediction. Various bucketing types have been proposed and described in detail in [3]. The Training bundle supports four types of bucketing: zero bucketing (i.e. fitting a single classifier), state-based bucketing, clustering-based bucketing and prefix length-based bucketing.

For each bucket of feature vectors, we train a predictive model using one of four supported machine learning techniques: decision tree, random forest, gradient boosting and extreme gradient boosting (XGBoost). For each technique, a user may manually enter the values of the most important hyperparameters. For example, when fitting a gradient boosting model, a user may choose the number of weak learners (trees), the number of features to be used for each split and the learning rate.

In order to accommodate users with varying degrees of expertise in machine learning and predictive process monitoring, the plugin bundle offers two training modes – basic and advanced. By default, the basic mode is activated wherein a user only needs to choose the log and prediction target. If the prediction target is based on the logi-

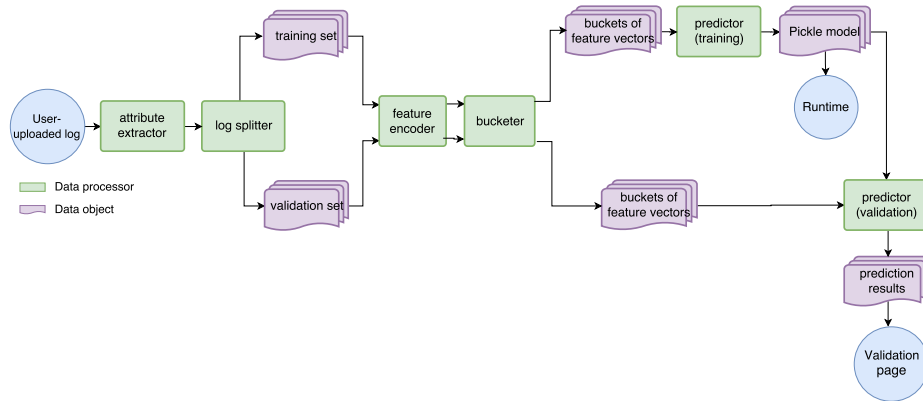


Fig. 4. High-level data flow diagram of the Training plugin bundle.

```

1 {
2   "case_id_col": "Line_ID",
3   "activity_col": "Activity_Name",
4   "timestamp_col": "Activity_End_Time",
5   "static_cat_cols": ["Supplier_ID", "Supplier_Location", "
   Delivery_Type", "Sector"],
6   "dynamic_cat_cols": ["Employee_ID", "weekday"],
7   "static_num_cols": ["Line_Total_Cost", "Target_Supply_Date_delta"]
8   "dynamic_num_cols": ["elapsed", "timesincelastevent", "open_cases"]
9   "ignore": ["Activity_Start_Time", "hour", "Milestone", "Key", "
   Target_Supply_Date", "Part_ID"],
10  "future_values": ["Late_Supply", "Delay_Rank"],
11 }

```

Fig. 5. Example training configuration file.

cal rule – whether the case duration will exceed the specified threshold, a user is also invited to key in the threshold value. For all the other settings – bucketing method, encoding method and prediction method and its hyperparameters – the default values which usually achieve the best prediction accuracy will be used. Experienced users may switch the advanced mode toggle and manually choose bucketing, encoding and prediction method settings or any sensible combination thereof. The latter is especially useful when a user wants to train and compare multiple models, e.g. using various sequence encoding methods.

The status of the trained models can be verified using the collapsible drawer in the right-hand corner. Upon the training completion, a serialized Python object in the pickle format is produced. It describes a trained predictive model and which includes:

- Configuration parameters of the predictors (whether it is a classifier or a regressor, what learning algorithm it uses)
- Definition of each column of the event log (static or dynamic, numeric or categorical). This information allows the Runtime plugin bundle to construct a feature vector from a given partial trace.
- The bucketing function, which given an input sample, allows us to determine which classifier/regressor should be used. Unless a zero bucketing is used, the pickle file has a different classifier/regressor per bucket.
- For each bucket, the trained model, ready to be taken as input by the selected classification algorithm, e.g. in the case of decision trees, the whole tree representation.

The predictive power of the trained model(s) can be evaluated on a held-out validation set. By default, a user will see the average accuracy across all partial traces after a certain number of events have completed. This evaluation method was also used in [1] and [3]. For classification tasks (e.g. prediction of Late Supply and Delay Rank), a user can choose which metrics to plot among accuracy score, F1 score and logarithmic loss. For regression tasks (e.g. remaining time), a user can choose between mean absolute error, either raw or normalized, and root mean square error. The accuracy of a particular model can be visually compared with that of other models trained for the same log and the same prediction target (Figure 6). Additionally, one can check a scatterplot of predicted vs. actual values (for regression tasks) or a confusion matrix (for classification tasks) and assess the relative importance of each feature for the chosen predictor.

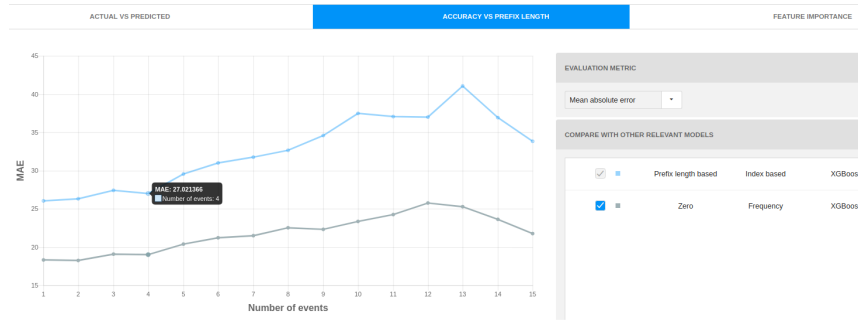


Fig. 6. Model validation page of the Training plugin bundle.

5 Runtime Plugin Bundle

Once the predictive models have been created, they can be deployed to the Runtime predictive monitoring environment of Apromore, to make predictions on ongoing cases. The Runtime plugin bundle can be used to stream an event log from the repository, or hook into an external stream. Either way, the input stream is transformed into a stream of predictions which is visualized in a Web-based dashboard. The transformation is implemented using the dataflow pipeline in Fig. 7.

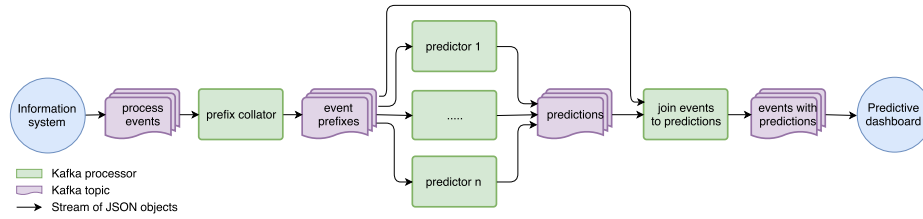


Fig. 7. High-level data flow diagram of the Runtime plugin bundle.

The pipeline is built on top of the open-source Apache Kafka stream processing platform. The “predictor” components of the pipeline are the predictive models from the Training plugin bundle. The “topic” components are network-accessible queues of JSON messages with publisher/subscriber support. This allows the computationally intense work of the predictors to be distributed across a cluster of networked computers, providing scalability and fault-tolerance. The “collator” component accumulates the sequence of events-to-date for each case, such that the prediction is a stateless function of the trained predictive model and of the case history. This statelessness is what allows the predictors to be freely duplicated and distributed. The “joiner” component composes the original events with the various predictions, ready for display on the dashboard.

The dashboard provides a list of both currently ongoing cases (colored in gray) as well as completed cases (colored in green), as shown in Fig. 8. For each case, it is also possible to visualize a range of summary statistics including the number of events in the case, its starting time and the time when the latest event in the case has occurred. For the ongoing cases, the Runtime plugin bundle provides the predicted values of the performance indicators the user wants to predict. For completed cases, instead, it shows

the actual values of the indicators. In addition to the table view, the dashboard offers other visualization options, such as pie charts for case outcomes and bar charts for case durations. It is also possible to export the predictions in a CSV file, for periodic reporting and for importing into third-party business intelligence tools.

Running cases

Completed cases

Completed events

Events per completed case

Average case duration

188

0

499

-

-

All events	Outcomes	Case duration	Remaining time	Case length								
Case	Events elapsed	Start time	Latest event time	Target supply date	Supplier Location	Delivery Type	Sector	Line Total Cost	Delay Rank	Late Supply	Next Activity	Predicted Completion
191315	4	2017-Aug-16 20:24	2017-Aug-22 19:54	2017-Nov-16	International	Sea	Manufacturing	\$2,048.00	91% Just in time	4%	91% Supply Date Recd	2017-Nov-09 21:15
191472	2	2017-Aug-20 23:53	2017-Aug-22 18:23	2017-Oct-28	International	Sea	Manufacturing	\$7,183.00	80% Just in time	34%	72% Order Confirmed	2017-Oct-26 07:25
182222	23	2017-Feb-10 01:27	2017-Aug-22 00:51	2017-Feb-12	International	Courier	Hi Tech	\$73,868.00	71% Source	100%	79% Goods Supplied	2017-Aug-27 18:08
190823	4	2017-Aug-09 23:38	2017-Aug-21 19:53	2017-Sep-19	International	Sea	Manufacturing	\$5,297.00	81% Just in time	26%	49% Delivered to Ship	2017-Sep-24 01:02
190822	5	2017-Aug-09 23:38	2017-Aug-21 19:53	2017-Sep-20	International	Sea	Manufacturing	\$3,558.00	92% Just in time	12%	47% Delivered to Ship	2017-Sep-28 15:16
190822	4	2017-Aug-09 23:38	2017-Aug-21 19:53	2017-Sep-20	International	Sea	Manufacturing	\$3,558.00	88% Just in time	14%	50% Delivered to Ship	2017-Sep-24 13:34
190873	2	2017-Aug-21 00:06	2017-Aug-21 00:06	2017-Nov-09	International	Sea	Hi Tech	\$883.00	88% Just in time	13%	75% Supply Date Recd	2017-Dec-06 18:38
190873	1	2017-Aug-21 00:06	2017-Aug-21 00:06	2017-Nov-09	International	Sea	Hi Tech	\$883.00	69% Just in time	18%	77% Order Confirmed	2017-Dec-06 09:08
190872	2	2017-Aug-21 00:06	2017-Aug-21 00:06	2017-Nov-03	International	Sea	Hi Tech	\$1,086.00	80% Just in time	27%	57% Order Confirmed	2017-Nov-10 05:42
190872	1	2017-Aug-21 00:06	2017-Aug-21 00:06	2017-Nov-03	International	Sea	Hi Tech	\$1,086.00	80% Just in time	32%	76% Order Confirmed	2017-Nov-23 11:49
190871	2	2017-Aug-21 00:06	2017-Aug-21 00:06	2017-Nov-01	International	Sea	Hi Tech	\$1,681.00	82% Just in time	22%	74% Supply Date Recd	2017-Nov-23 20:35
190871	1	2017-Aug-21 00:06	2017-Aug-21 00:06	2017-Nov-01	International	Sea	Hi Tech	\$1,681.00	84% Just in time	27%	77% Order Confirmed	2017-Nov-24 10:48
190870	2	2017-Aug-21 00:06	2017-Aug-21 00:06	2017-Oct-29	International	Sea	Hi Tech	\$2,162.00	62% Just in time	30%	75% Supply Date Recd	2017-Nov-21 07:21
190870	1	2017-Aug-21 00:06	2017-Aug-21 00:06	2017-Oct-29	International	Sea	Hi Tech	\$2,162.00	63% Just in time	35%	78% Order Confirmed	2017-Nov-23 17:25
190869	2	2017-Aug-21 00:06	2017-Aug-21 00:06	2017-Nov-04	International	Sea	Hi Tech	\$3,957.00	86% Just in time	20%	73% Supply Date Recd	2017-Dec-03 18:22

<<

<

1

/

25

>

>>

1 - 20 / 498

Fig. 8. Main view of the dashboard in the Runtime plugin bundle.

Process workers and operational managers can set some process performance targets and subscribe to a stream of warnings and alerts generated whenever these targets are predicted to be violated. Thus, they will be capable of making informed, data-driven decisions to get a better control of the process executions. This is especially beneficial for processes where process participants have more leeway to make corrective actions (for example, in a lead management process).

6 Conclusion

Through the integration with Nirdizati, Apromore offers a configurable full-stack Web tool that supports users in selecting and tuning various prediction models, and that enables the continuous prediction of different process performance indicators at runtime. Predictions can be presented visually in a dashboard or exported for periodic reporting.

Video demos of the model training and of the runtime functionality can be found at <http://youtu.be/CYwLTKAbF1Q> and at <http://youtu.be/Q4WVebqJzUI>. The source code is available under the LGPL version 3.0 license at <https://github.com/apromore/ApromoreCode>.

References

1. Leontjeva, A., Conforti, R., Francescomarino, C.D., Dumas, M., Maggi, F.M.: Complex symbolic sequence encodings for predictive monitoring of business processes. BPM (2015)
2. Márquez-Chamorro, A.E., Resinas, M., Ruiz-Cortes, A.: Predictive monitoring of business processes: a survey. IEEE TTSC (2017)
3. Teinmaa, I., Dumas, M., La Rosa, M., Maggi, F.M.: Outcome-oriented predictive process monitoring: Review and benchmark. CoRR abs/1707.06766 (2017)
4. Verenich, I.: A general framework for predictive business process monitoring. CAiSE Doctoral Consortium (2016)