

Object-centric Process Predictive Analytics

Riccardo Galanti^{a,b}, Massimiliano de Leoni^b, Nicolò Navarin^b, Alan Marazzi^a

^aIBM, Bologna, Italy

^bUniversity of Padua, Padua, Italy

Abstract

Object-centric processes (a.k.a. Artifact-centric processes) are implementations of a paradigm where an instance of one process is not executed in isolation but interacts with other instances of the same or other processes. Interactions take place through bridging events where instances exchange data. Object-centric processes are recently gaining popularity in academia and industry, because their nature is observed in many application scenarios. This poses significant challenges in predictive analytics due to the complex intricacy of the process instances that relate to each other via many-to-many associations. Existing research is unable to directly exploit the benefits of these interactions, thus limiting the prediction quality. This paper proposes an approach to incorporate the information about the object interactions into the predictive models. The approach is assessed on real-life object-centric process event data, using different KPIs. The results are compared with a naïve approach that overlooks the object interactions, thus illustrating the benefits of their use on the prediction quality.

Keywords: Predictive Analytics, Object-centric Process, Gradient Boosting, Artifact-centric Process, Process Mining

1. Introduction

Traditionally, processes are seen as being instantiated in cases that are constituted by single flows that are executed in isolation. However, inter-organization processes

Email addresses: riccardo.galanti@ibm.com (Riccardo Galanti), deleoni@math.unipd.it (Massimiliano de Leoni), nnavarin@math.unipd.it (Nicolò Navarin), alan.marazzi@ibm.com (Alan Marazzi)

are oftentimes more complex: usually, several instances of different processes are being executed at the same time and they may interact with each other. In fact, the situation is more similar to choreographies where one instance of a process P_1 interacts and synchronizes with several instances of a second process P_2 , and the other way around: one instance of P_2 might synchronize with multiple instances of P_1 . The situation can be even more complex: instances of P_2 may in turn interact with instances of some P_3 , and so on. For instance, consider a retail shop in Padua (Italy): several customers may order products manufactured in a factory in Brisbane (Australia). The factory associates many customer orders to a single manufacturer order to save money. Also, the same customer orders can include products from different manufacturers in different parts of the globe. Customer's and manufacturer's orders are managed via instances of different processes: one instance of customer-order process can be associated to several of manufacturer-order process, and the other way round: each manufacturer-order process instance may be associated to many consumer-order ones.

The paradigm of object-centric processes (a.k.a. artifact-centric processes) is gaining more and more momentum in the recent years to model inter-organizational processes more naturally [1, 2, 3]. Any process execution materializes itself as a set of instances of the same/different processes that represent the life cycles of different objects (a.k.a. artifacts) that contribute to the process execution (e.g., the order and the delivery object). These processes for the different objects run independently and synchronize through some bridging events to exchange data needed to progress further. The importance of object-centric process approaches is particularly evident in the domain of Process Mining: the IEEE Task Force has recently reported on the results of a survey with academics, practitioners, consultants and vendors in Process Mining, according to which only 33% of the respondents consider forcing to pick one single case identifier to be a minor problem [4].

This paper focuses on object-centric process predictive analytics. This is a challenging problem that cannot be tackled via the current research in process predictive analytics, since the latter relies on the assumption that the instances refer to a single process. In a nutshell, the main idea is that a main object type (e.g., *Customer Order*) is chosen as viewpoint. The complex object interaction is unfolded in traces around

the viewpoint: one trace is created per object o of the viewpoint type, and includes the events related to that object (i.e., process instance), and related to the objects of the same or different types that synchronize with o , directly or indirectly. When the complex interaction is unfolded in a multiset of traces, we can specialize the current state of the art in predictive analytics.

Process stakeholders define the Key Performance Indicators (KPIs) using domain knowledge: this corresponds to instantiating a function $\mathcal{T}(\sigma)$ that returns the KPI value for any trace σ . Since the viewpoint determines how events are grouped in traces, the determination of the appropriate KPI influences which viewpoint to choose. As an example, if the KPI relates to the duration of an order process, the viewpoint needs to be the order object.

The proposed approach is assessed through experiments on a real event log extracted from an object-centric process executed by an utility provider company in Italy. Experiments were conducted on different KPIs, illustrating that the complex interactions of object-centric processes need to be taken into account when predicting, as this allows to consistently improve the predictive performances over simpler techniques.

Section 2 introduces the preliminary concepts that are used later on: object-oriented and single-id event logs, KPIs, and the traditional problem of process predictive analytics. Section 3 illustrates the procedure to create single-id event logs from object-oriented event logs, also illustrating how aggregation attributes can be additionally used to better capture the object-to-object synchronization dynamics. Section 4 sketches some implementation details and extensively reports on the experimental phase, while related works are discussed in Section 5. Finally, Section 6 concludes the paper, summarizing the paper context and contribution, and delineating some potential avenues of future work.

2. Preliminaries

2.1. Object-Centric Event Logs

Object-centric processes are carried on with the support of one or more information systems. It is possible to extract the history of past executions into a transactional data

set organized in form of object-centric event logs [5]:

Definition 2.1 (Object-Centric Event Log). *Let T be the universe of the timestamps.*

An object-centric event log is a tuple $L = (E, A, AN, AV, AT, OT, O, \pi_{typ}, \pi_{act}, \pi_{time}, \pi_{vmap}, \pi_{omap}, \pi_{otyp}, \pi_{ovmap}, <)$ such that:

- E is the set of event identifiers,
- A is the set of activity names,
- AN is the set of attributes names,
- AV is the set of attribute values (with the requirement that $AN \cap AV = \emptyset$),
- AT is the set of attribute types,
- OT is the set of object types,
- O is the set of object identifiers,
- $\pi_{typ} : AN \cup AV \rightarrow AT$ is the function associating an attribute name or value to its corresponding type,
- $\pi_{act} : E \rightarrow A$ is the function associating an event identifier to its activity,
- $\pi_{time} : E \rightarrow T$ is the function associating timestamps to event identifiers,
- $\pi_{vmap} : E \rightarrow (AN \not\rightarrow AV)$ is the function associating every event identifier $e \in E$ to a variable-to-value assignment function val such that, for each attribute $a \in AN$ in the domain of val , $val(a)$ indicates the value assigned to a by e ,¹
- $\pi_{omap} : E \rightarrow 2^O$ is the function associating an event identifier to a set of related object identifiers,
- $\pi_{otyp} : O \rightarrow OT$ assigns precisely one object type to each object identifier,
- $< \subseteq (E \times E)$ is a partial order of events.²

¹The notation $\not\rightarrow$ indicates a partial function.

²Typically, the partial order is induced by the timestamp, i.e., $e' < e'' \iff \pi_{time}(e') < \pi_{time}(e'')$.

However, we do not require to make that assumption.

Table 1: Example of object-centric event log. Each row represents an event. The blank spaces represent attributes missing values.

id	activity	timestamp	Contract	Requisition	Order	Receipt	Invoice	user	order_price	order_delivery_month	order_purch_group	rec_quantity
e1	Contract Line Creation	2017-07-11 9:00	c1					CO01				
e2	Purch Contract Item Material Group Changed	2017-07-14 11:00	c1					CO01				
e3	Purchase Requisition Line Created	2017-07-15 12:00	c1	rq1				A456				
e4	Purchase Requisition Line Created	2017-07-15 15:00	c1	rq2				A457				
e5	Purchase Order Line Creation	2017-07-16 15:00	c1		o1			A458	100	7	100.L50	
e6	Purchase Order Line Creation	2017-07-17 15:00		rq1	o2			A458	200	8	100.L51	
e7	Purchase Order Line Creation	2017-07-18 15:00		rq2	o3			A458	300	8	100.L52	
e8	Goods Line Registered	2017-07-22 15:00			o1	r1		A456	100	7	100.L50	10
e9	Invoice Receipt	2017-07-22 16:00					i1	A125				
e10	Purchase Requisition Group Changed	2017-07-22 19:00		rq1				A456				
e11	Purchase Order Line Creation	2017-07-23 9:00		rq1	o4			A458	600	8	100.L51	
e12	Goods Line Registered	2017-07-23 15:00			o2	r2		A456	100	8	100.L50	10
e13	Invoice Registered	2017-07-29 11:00				r1,r2	i1	A125				10
e14	Invoice Cleared	2017-07-30 12:00					i1	A125				
e15	Goods Line Registered	2017-07-31 15:00			o4	r3		A456	600	8	100.L51	10
e16	Invoice Registered	2017-08-10 11:00				r2,r3	i2	A125				10
e17	Invoice Cleared	2017-08-15 14:00					i2	A125				
e18	Goods Line Registered	2017-08-16 15:00			o3	r4		A456	300	8	100.L52	5
e19	Purchase Requisition Supplier Changed	2017-08-16 17:00		rq2				A456				
e20	Invoice Registered	2017-08-18 11:00				r4	i3	A125				5
e21	Invoice Cleared	2017-08-20 14:00					i3	A125				

Example 2.1. Table 1 shows an excerpt of an object-centric event log of an Italian utility provider company. It consists of five object types, each with its own object identifier: *Contract*, *Requisition*, *Order*, *Receipt*, *Invoice*. The first is *Contract*, which is the process concerning the stipulation of a contract with a customer, possibly followed by a *Requisition*, which is an optional process executed when the order needs a purchase requisition. The *Order* process consists of several activities representing mainly quantity, price, or date modifications of the order, eventually approved by the Head of the department. The *Receipt* process is then related to the receiving of the goods or the services requested, followed by the *Invoice* process, which includes everything related to payments. Some events are associated to a single object identifier, others have multiple (i.e., the so-called bridge events) that enable the synchronization and data exchange between objects. Figure 1 illustrates how objects are related to each other for synchronization and data exchanges. Note that relationships can be of many-to-many or many-to-one nature. Events are associated with attributes, features based on the properties of requisitions, orders (e.g., *order_price*), receipts (e.g., *receipt_quantity*), and invoices.

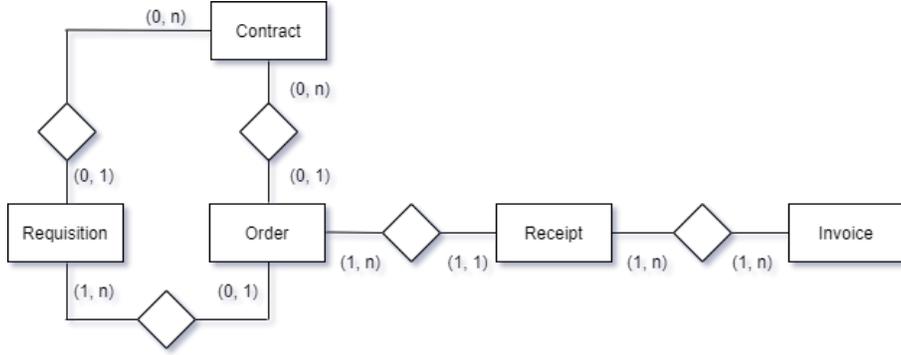


Figure 1: Diagram representing cardinality between the different object types in the considered object-centric event log. For each object type, the cardinality with the subsequent or the previous object type is represented as $(min_cardinality, max_cardinality)$

2.2. Single-Id Process Predictive Analytics

In this section, we discuss the typical techniques adopted to train a predictive model starting from an event log with a single identifier.

Definition 2.2 (Single-id Event Log). A single-id event log is a tuple: $\mathcal{L} = (E, T, A, AN, AV, AT, \pi_{typ}, \pi_{act}, \pi_{time}, \pi_{vmap})$ consisting of a set E of event identifiers, a set $T \subset E^*$ of traces, i.e. sequences of event identifiers, a set A of activity names, sets AN and AV of attribute names and values, a function π_{typ} associating attribute names and values to types, $\pi_{act} : E \rightarrow A, \pi_{time} : E \rightarrow T, \pi_{vmap} : E \rightarrow (AN \not\rightarrow AV)^3$ associating each event identifier to event's activity, timestamp, and attribute assignment, respectively.

The aim is to predict the value of a key performance indicator (KPI), which depends on the specific process domain. This corresponds to providing the definition a KPI function.

Definition 2.3 (KPI Function). Let $\mathcal{L} = (E, T, A, AN, AV, AT, \pi_{typ}, \pi_{act}, \pi_{time}, \pi_{vmap})$ be a single-id event log. Let \mathcal{W}_K be the set of possible KPI values. A KPI

³Notation $A \not\rightarrow B$ is used to highlight a partial function from some elements of a domain A to elements of B

is a function $\mathcal{T}_{\mathcal{L}} : E^* \times \mathbb{N} \dashrightarrow \mathcal{W}_K$ such that, given a trace $\sigma \in E^*$ and an integer index $i \leq |\sigma|$, $\mathcal{T}_{\mathcal{L}}(\sigma, i)$ returns the KPI value of σ after the occurrence of the first i events.⁴

Note that our KPI definition assumes it to be computed a posteriori, when the execution is completed and leaves a complete trail as a certain trace σ . In many cases, the KPI value is updated after each activity execution, which is recorded as next event in trace; however, other times, this is only known after the completion. In the remainder, when clear from the context, we often omit the subscript \mathcal{L} .

We aim to be generic and account for all relevant domains. Given a trace $\sigma = \langle e_1, \dots, e_n \rangle$ that records a complete process execution, the following are example of three potential KPI definitions:

Remaining Time. This corresponds to the situation in which, after executing a sequence of i events, the KPI measures how long the process execution will still last to completion: $\mathcal{T}_{remaining}(\sigma, i) = \pi_{time}(e_n) - \pi_{time}(e_i)$, namely the difference between the timestamp of latest future trace event e_n and that of the last occurred event e_i .

Activity Occurrence. It measures whether a certain activity is going to eventually occur in the future, such as an activity *Open Loan* in a loan-application process. The corresponding KPI definition for the occurrence of an activity A is $\mathcal{T}_{occur_A}(\sigma, i)$, which is equal to true if activity A occurs in $\langle e_{i+1}, \dots, e_n \rangle$ and $i < n$; otherwise false.

Customer Satisfaction. This is a typical KPI for several service providers. Let us assume, without losing generality, to have a trace $\sigma = \langle e_1, \dots, e_n \rangle$ where the satisfaction is known at the end, e.g. through a questionnaire. Assuming the satisfaction level is recorded with the last event - say $e_n(sat)$. Then, $\mathcal{T}_{cust_satisf}(\sigma, i) = e_n(sat)$.

We can define the prediction problem on unfolded logs as follows.

⁴Given a sequence X , $|X|$ indicates the length of X .

Definition 2.4 (Prediction Problem on Single-id Event Logs). *Let \mathcal{L} be a single-id event log that records the execution of a given process, for which a KPI $\mathcal{T}_{\mathcal{L}}$ is defined. Let $\sigma = \langle e_1, \dots, e_k \rangle$ be the trace of a running case, which eventually will complete as $\sigma_T = \langle e_1, \dots, e_k, e_{k+1} \dots, e_n \rangle$. The prediction problem can be formulated as forecasting the value of $\mathcal{T}_{\mathcal{L}}(\sigma_T, i)$ for all $k < i \leq n$.*

In the process mining literature, this problem has been faced with different machine learning models [6, 7, 8, 9, 10, 11].

The training set is composed by pairs $(x, y) \in \mathcal{X} \times \mathcal{Y}$ where \mathcal{X} encodes the independent variables (also known as **features**) with their values, and \mathcal{Y} is the the dependent variable with its value (i.e. the value to predict). Process predictive analytics requires a KPI definition $\bar{\mathcal{T}}$ as input (cf. Definition 2.3). Let $\mathcal{W}_K = \text{img}(\bar{\mathcal{T}})$ be the domain of possible KPI values (i.e. the image/codomain of $\bar{\mathcal{T}}: \mathcal{Y} = \mathcal{W}_K$. Afterwards, each prediction technique requires the definition of the domain \mathcal{X} and a **trace-to-instance encoding function** $\rho: E^* \rightarrow \mathcal{X}$, which maps each (prefix of a) trace σ in an element $\rho(\sigma) \in \mathcal{X}$.

The prediction model is trained off-line via a dataset \mathcal{D} that is created from an event log \mathcal{L} as follows. Each prefix σ of each each trace $\sigma_T \in \mathcal{L}$ generates one distinct item in \mathcal{D} consisting of a pair $(x, y) \in (\mathcal{X} \times \mathcal{Y})$ where $x = \rho(\sigma)$ and $y = \bar{\mathcal{T}}(\sigma_T, |\sigma|)$. Once the dataset item of every trace prefix is created, the model is trained. The resulting prediction model (a.k.a. predictor) can be abstracted as an oracle function $\Phi_{\mathcal{D}}: \mathcal{X} \rightarrow \mathcal{Y}$.

This paper does not focus on comparing different machine learning techniques. Instead, we focus on how it is possible to leverage on methodologies developed for single-process prediction problems to tackle the object-centric process prediction problem. In the implementation and experiments, we use the Catboost method [12], a high-performance open source framework for gradient boosting on decision trees, but different types of predictive models could be learnt. The choice fell onto Catboost because experiments show that it outperforms and solves limitations of current state-of-the-art implementations of gradient boosted decision trees [12]. It is backed by solid theoretical results that explain how strong predictors can be built by iteratively com-

binning weaker models (base predictors) in a greedy manner. Catboost, in particular, at each iteration t of the algorithm performs a random permutation of the features and a tree is constructed on the basis of it. Moreover, for each split of a tree, CatBoost combines (concatenates) all categorical features (and their combinations) already used for previous splits in the current tree with all categorical features in the dataset.

In the domain of Catboost learning, the definition of the trace-to-instance encoding function requires the intermediate concept of *event-to-tuple function* $\zeta_{\mathcal{L}} : E \rightarrow A \times AT \times (AV)^w$, which encodes each event of a single-id event log $\mathcal{L} = (E, T, A, AN, AV, AT, \pi_{typ}, \pi_{act}, \pi_{time}, \pi_{vmap})$, where $w = |AN|$ is the number of attributes defined in the event log \mathcal{L} . The event-to-tuple function is defined as follows. In particular, indicated the concatenation of two tuples with \oplus , given an event in \mathcal{L} with identifier e_i , $\zeta_{\mathcal{L}}(e) = [\pi_{act}(e), \pi_{time}(e)] \bigoplus_{v \in AN} \pi_{vmap}(v)$.⁵

In Catboost, the trace-to-instance encoding function also considers the history of each partial trace σ by considering the number of times that each activity has been performed in σ . Consequently, we define the function $\rho_{\mathcal{L}}^{agg}(\langle e_1, \dots, e_m \rangle)$; here, for each activity $a \in A$, one dimension exists in $\rho_{\mathcal{L}}^{agg}(\sigma) : E^* \rightarrow (\mathbb{N})^{|A|}$ that takes on a value equal to the number of events $\bar{e}\sigma$ that refer to a , i.e. such that $\pi_{act}\bar{e}$. The function $\rho_{\mathcal{L}}$ is then defined as: $\rho_{\mathcal{L}}(\langle e_1, \dots, e_m \rangle) = \rho_{\mathcal{L}}^{agg}(\langle e_1, \dots, e_m \rangle) \oplus \zeta_{\mathcal{L}}(e_m)$.

2.3. Explanations of Process Predictions

Several prediction models, including Catboost, are black boxes, making it difficult to explain the predictions, namely to determine the degree with which each feature influences the predictions. Explaining the predictions is beneficial to build user trust in the prediction model. The remainder briefly summarizes the basic concepts behind the framework for prediction explanation that has been introduced in [13]. This framework has been used during the evaluation to illustrate that several features related to the object interaction have a significant impact on the predictions, thus consequently showing their relevance to improve the prediction accuracy.

⁵To keep the explanation simple, we assume that the enumerations of all variables v in AN are always returned consistently, as if there is a total order among the variables (e.g., the alphabetical order).

The framework in [13] leverages on the Shapley Values [14], which is a game theory approach to fairly distribute the payout among the players that have collaborated in a cooperative game. The assumption is that the features from an instance correspond to the players, and the payout is the difference between the prediction made by the predictive model and the average prediction (also called *base value*). Intuitively, given a predicted instance, the Shapley Value of a feature expresses how much the feature value contributes to the model prediction [15]:

Definition 2.5 (Shapley Value). *Let $X = \{x_1, \dots, x_n\}$ be a set of features. The Shapley value for feature x_i is defined as:*

$$\psi_i = \sum_{S \subseteq \{x_1, \dots, x_m\} \setminus \{x_i\}} \frac{|S|!(p-|S|-1)!}{p!} (val(S \cup \{x_i\}) - val(S))$$

where $val(X')$ is the so-called payout for only using the set of feature values in $X' \subset X$ in making the prediction.

Intuitively, the formula in Definition 2.5 evaluates the effect of incorporating the feature value x_i into any possible subset of the feature values considered for prediction. In the equation, variable S runs over all possible subsets of feature values, the term $val(S \cup \{x_i\}) - val(S)$ corresponds to the marginal value of adding x_i in the prediction using only the set of feature values in S , and the term $\frac{|S|!(p-|S|-1)!}{p!}$ corresponds to all the possible permutations with subset size $|S|$, to weight different sets differently in the formula. This way, all possible subsets of attributes are considered, and the corresponding effect is used to compute the Shapley Value of x_i .

The starting point for the explainable framework is the trace-to-instance encoding function $\rho : E^* \rightarrow \mathcal{X}$ (cf. Section 2.2), and a single-id event $\log \mathcal{L} = (E, T, A, AN, AV, AT, \pi_{typ}, \pi_{act}, \pi_{time}, \pi_{vmap})$.

Let us recall that given a trace $\sigma = \langle e_1, \dots, e_m \rangle \in T$, $\rho(\sigma) = [x^1, \dots, x^n]$, each feature f^i has an associated value x^i . As mentioned in Section 2.2, a feature f^i can be of different nature, such as a process attribute, a timestamp, or the number of executions of an activity in σ . The prediction model is built over the multiset $\uplus_{\sigma' \in T} \rho(\sigma')$.

When applied for prediction explanations, the Shapley values for a trace σ are computed over tuple $\rho(\sigma) = [x^1, \dots, x^n]$, thus resulting in a tuple of Shapley values

$\Psi = [\psi^1, \dots, \psi^n]$, with ψ^i being the Shapley value of feature f^i . In accordance with the Shapley values theory, the explanation of ψ^i is as follows: since feature $f^i = x^i$, the KPI prediction deviates ψ^i units from the average KPI value of the event-log traces.

The computation of the Shapley value is repeated for each trace of \mathcal{L} . However, if f^i is numerical, several different values can be observed for f^i , yielding a large number of explanations $f^i = x_1^i, \dots, f^i = x_k^i$. Some of these explanations are equivalent from a domain viewpoint: e.g., $amount = 10000$, $amount = 10050$ might be referring to the same class of amount in a loan application. Therefore, q representative values w_1^i, \dots, w_q^i are selected out of values x_1^i, \dots, x_k^i (namely, with $q \ll k$) so as to obtain explanations of type $f^i < w_1^i, w_1^i \leq f^i < w_2^i, \dots, f^i \geq w_q^i$. Values w_1^i, \dots, w_q^i can be obtained taking the boundaries of the buckets obtained via discretization techniques. In particular, our implementation operationalizes a discretization of each feature f^i on the basis of decision/regression as follows. The training set consists of tuple with only two features: f^i used as the independent variable, and the KPI as target/dependent variable. The values observed at the splits of the tree nodes induce the boundaries and, consequently, the buckets.

While an exact computation of the Shapley values requires to consider all combinations of features, efficient estimations can be obtained through polynomial algorithms that use greedy approaches [15].

The discussion so far focused on the Shapley values for single (prefixes of) traces. It is possible to show the distribution over multiple traces via, e.g., boxplots. Examples will be shown in Section 4 (cf. Figures 2, 3 and 4).

3. Predictive Analytics in Object-Centric Processes

The starting point is an object-centric log $L = (E, A, AN, AV, AT, OT, O, \pi_{typ}, \pi_{act}, \pi_{time}, \pi_{vmap}, \pi_{omap}, \pi_{otyp}, \pi_{ovmap}, <)$. Our object-centric process prediction requires analysts to decide a so-called *viewpoint*, which is an object type $o_t \in OT$ of the process (e.g., *Requisition*). This defines how the events in an object-centric event log are aggregated to form traces of a single-id event log.

A single-id event log $\mathcal{L} = (E', T, A', AN', AV', AT', \pi'_{typ}, \pi'_{act}, \pi'_{time}, \pi'_{vmap})$ is

created from a chosen viewpoint $o_t \in OT$ as follows. The trace set T contains one trace σ_o for each object $o \in O$ such that $\pi_{otyp}(o) = o_t$. To determine which events to include in σ_o , we compute the timestamp of the first event in E that has o as one of the object identifiers:

$$t_o = \min_{e \in E. o \in \pi_{omap}(e)} \pi_{time}(e) \quad (1)$$

Trace σ_o will include every event $e \in E$ with timestamp larger than or equal to t_o such that it at least contains o as identifier (namely, such that $\{o\} \subseteq \pi_{omap}(e)$) or contains an identifier of an object o' in a certain set $R_L^+(o)$ of related objects (namely such that $R_L^+(o) \cap \pi_{omap}(e) \neq \emptyset$). The constraint to exclude events with timestamp smaller than t_o is motivated by the fact that, if an event \bar{e} precedes the first event e_1 of σ_o , \bar{e} is unrelated to the behavior in σ_o : if \bar{e} had influenced or was influenced by the behavior in σ_o , \bar{e} would have occurred after e_1 , which is, in fact, intended as the creation event of the process related to life cycle of o .

The set $R_L^+(o)$ is constructed as follows. Let us define $R_{L,o_t}^1(o)$ as the set of objects directly related to o of type different than o_t :

$$R_{L,o_t}^1(o) = \{o' \in O : \exists e \in E. \{o, o'\} \subseteq \pi_{omap}(e) \wedge \pi_{otyp}(o') \neq o_t\}.$$

We can also define the set of objects of type different than o_t that are indirectly related to o via a “bridge” object o' of type different than o_t :

$$R_{L,o_t}^2(o) = \bigcup_{o' \in R_{L,o_t}^1(o)} R_{L,o_t}^1(o') \quad (2)$$

The definition above can also be extended for any $R_{L,o_t}^i(o)$ with $i > 1$ as follows:

$$R_{L,o_t}^i(o) = \bigcup_{o' \in R_{L,o_t}^{i-1}(o)} R_{L,o_t}^1(o') \quad (3)$$

The set $R_L^+(o)$ can thus be defined as the union, for all integer indexes $i \geq 1$, of the sets of the objects of type different than $\pi_{otyp}(o)$ that are related to o via $(i - 1)$ bridging events of type different than $\pi_{otyp}(o)$:

$$R_L^+(o) = \bigcup_{i=1}^{\infty} R_{L,\pi_{otyp}(o)}^i(o) \quad (4)$$

Table 2: Unfolded event log obtained from the object-centric event log in Table 1 with the proposed object-centric approach when considering the Requisition viewpoint. The horizontal lines split the different traces. The blank spaces represent attributes missing values.

id	activity	timestamp	Contract	Requisition	Order	Receipt	Invoice	user	order_price	order_delivery_month	order_purch_group	rec_quantity
e3	Purchase Requisition Line Created	2017-07-15 12:00	c1	rq1				A456				
e4	Purchase Requisition Line Created	2017-07-15 15:00	c1	rq2				A457				
e5	Purchase Order Line Creation	2017-07-16 15:00	c1		o1			A458	100	7	100.L50	
e6	Purchase Order Line Creation	2017-07-17 15:00		rq1	o2			A458	200	8	100.L51	
e7	Goods Line Registered	2017-07-22 15:00			o1	r1		A456	100	7	100.L50	10
e9	Invoice Receipt	2017-07-22 16:00					ii	A125				
e10	Purchase Requisition Group Changed	2017-07-22 19:00		rq1				A456				
e11	Purchase Order Line Creation	2017-07-23 9:00		rq1	o4			A458	600	8	100.L51	
e12	Goods Line Registered	2017-07-23 15:00			o2	r2		A456	200	8	100.L50	10
e13	Invoice Registered	2017-07-29 11:00				r1,r2	ii	A125				10
e14	Invoice Cleared	2017-07-30 12:00					ii	A125				
e15	Goods Line Registered	2017-07-31 15:00			o4	r3		A456	600	8	100.L51	10
e16	Invoice Registered	2017-08-10 11:00				r2,r3	i2	A125				10
e17	Invoice Cleared	2017-08-15 14:00					i2	A125				
e4	Purchase Requisition Line Created	2017-07-15 15:00	c1	rq2				A457				
e7	Purchase Order Line Creation	2017-07-18 15:00		rq2	o3			A458	300	8	100.L52	
e18	Goods Line Registered	2017-08-16 15:00			o3	r4		A456	300	8	100.L52	5
e19	Purchase Requisition Supplier Changed	2017-08-16 17:00		rq2				A456				
e20	Invoice Registered	2017-08-18 11:00				r4	i3	A125				5
e21	Invoice Cleared	2017-08-20 14:00					i3	A125				

Table 3: Unfolded event log obtained when considering the Requisition viewpoint from the object-centric event log in Table 1 with the existing approach, assuming each instance belongs to a single process. The horizontal lines split the different traces. The blank spaces represent attributes missing values.

id	activity	timestamp	Contract	Requisition	Order	Receipt	Invoice	user	order_price	order_delivery_month	order_purch_group	rec_quantity
e3	Purchase Requisition Line Created	2017-07-15 12:00	c1	rq1				A456				
e6	Purchase Order Line Creation	2017-07-17 15:00		rq1	o2			A458	200	8	100.L51	
e10	Purchase Requisition Group Changed	2017-07-22 19:00		rq1				A456				
e11	Purchase Order Line Creation	2017-07-23 9:00		rq1	o4			A458	600	8	100.L51	
e4	Purchase Requisition Line Created	2017-07-15 15:00	c1	rq2				A457				
e7	Purchase Order Line Creation	2017-07-18 15:00		rq2	o3			A458	300	8	100.L52	
e19	Purchase Requisition Supplier Changed	2017-08-16 17:00		rq2				A456				

The application of the aforementioned procedure to each object o of the viewpoint type creates the set T of traces, which in turn induces the the other elements of the \mathcal{L} tuple as follows. The set $E' = \cup_{\sigma_o \in T} \cup e' \in \sigma_o e$ of event identifier contains the event identifiers in T . The domain of attribute names, values and types is the same as in the object-centric log: $AN' = AN$, $AV' = AV$, $AT' = AT$, and $\pi'_{typ} = \pi_{typ}$. Functions $\pi'_{act}, \pi'_{time}, \pi'_{vmap}$ are restricted over domain E' , namely for each $e' \in E'$ $\pi'_{act}(e) = \pi_{act}(e)$, $\pi'_{time}(e) = \pi_{time}(e)$, $\pi'_{vmap}(e) = \pi_{vmap}(e)$.

Example 3.1. *Let us assume to unfold the object-oriented event log in Table 1, using Requisition as viewpoint. The result is presented in Table 2. Traces are split through horizontal lines. There are two traces, as many as the number of requisitions. In particular, since in the first trace we focused on the requisition with identifier rq1,*

we first considered the events containing $rq1$ as identifier ($e3, e6, e10, e11$); afterwards, we considered the events related transitively to $rq1$, i.e. with at least one identifier of an object in $R_L^+(rq1)$. For instance, we considered event $e5$ because it contains the contract identifier $c1$ that is related to $rq1$ via event $e3$, which contains both $rq1$ and $c1$ among the identifiers. Please note that, e.g., in the first trace we do not consider event $e19$ because $e19$ only contains the object identifier $rq2$, which is of the same type as $rq1$, namely the requisition object type. Neither do we consider such events as $e18$ because, e.g., $e18$ is only associated to $rq1$ via a transitive relation that goes through $rq2$, which is of the same type as $rq1$. Event $e4$ is instead included in the first trace because it is associated to the contract $c1$, which is associated transitively to the requisition $rq1$. Event $e4$ is also in the second trace because it is associated to $rq2$, since $rq2$ is within the object identifiers of $e4$. It is worth noting that events $e1$ and $e2$ are excluded from the first and second trace since their timestamps are smaller than $e3$ or $e4$, which are respectively the first events with $rq1$ or $rq2$ as object identifiers.

Conversely, existing techniques would be unable to deal with multiple object types, namely with multiple interleaving processes; they would only restrict to one single process, which would likely be the process related to objects of type *Requisition*. Therefore, there would again be two traces, one per requisition. Differently from our approach, the event log would only contain the events that present the first or the second requisition among the identifiers. The resulting log is in Table 3.

Summarising, the approach starts from an object-centric event log $L = (E, A, AN, AV, AT, OT, O, \pi_{typ}, \pi_{act}, \pi_{time}, \pi_{vmap}, \pi_{omap}, \pi_{otyp}, \pi_{ovmap}, <)$, which is unfolded to a single-id event log $\mathcal{L} = (E', T, A', AN', AV', AT', \pi'_{typ}, \pi'_{act}, \pi'_{time}, \pi'_{vmap})$ around a viewpoint, namely an object type $o_t \in OT$ of the process (e.g., *Requisition*). The viewpoint defines how the events in an object-centric event log are aggregated to form traces of a single-id event log.

These traces can be used to train and test a proper prediction model using the off-the-shelf techniques discussed in Section 2.2. However, this unfolding to a single-id event log does not preserve information about the number of objects correlated to the

viewpoint object. In fact, it excludes the information about the attributes associated with correlated objects, and their respective value. For instance, let us assume that we aim to use the total processing time of requisitions as KPI, and this time is somehow correlated to the number of orders associated to the requisitions (e.g., a requisition for more orders takes longer to be processed). The prediction of this KPI could be more accurate if the number of orders objects associated with the requisition were used to train and use the prediction model.

To mitigate this information loss, we extend our first approach to a richer one, where we introduce aggregation attributes that synthesize additional interaction information. Given a prefix σ'_o , of the trace σ_o for a viewpoint object o , the following aggregation attributes are included.

1. A given attribute a that can take different values v_1, \dots, v_n in a set $\{o'_1, \dots, o'_m\} \subseteq R_L^+(o)$ of objects related to o . The standard encoding would only retain one value for a , namely the value observed in the latest event in σ_o that assigns a value to a . We define:
 - if a is numerical (i.e., $\pi_{typ}(a) \subseteq \mathbb{R}$), one feature f is added to the feature set that contains an aggregated value $\psi(v_1, \dots, v_n)$ summarizing observed values. The aggregation function $\psi(\vec{v})$ may be customized depending on the domain, such as the average value in \vec{v} ;
 - if a is categorical, a feature is added for each pair (a, v_i) with value defined as the ratio between the number of attribute assignments to the value v_i over the total assignments of a .
2. For each object type $o'_t \in OT$, one feature encodes the number of objects of type o'_t associated with events of σ'_o , namely $|\{o' \in O : \pi_{otyp}(o) = o'_t \wedge \exists e \in \sigma'_o. o' \in \pi_{omap}(e)\}|$.
3. For each object type $o'_t \in OT$ and for each activity $a \in A$ one feature is added with value equal to the percentage of objects of type o'_t for which activity a has occurred at least once in trace σ'_o .

Table 4: Unfolded event log obtained from the object-centric event log in Table 1 with the proposed object-centric approach when considering the Requisition viewpoint and the aggregation attributes. The horizontal lines split the different traces. The blank spaces represent attributes missing values.

id	activity	timestamp	Contract	Requisition	Order	Order- delivery- month	Order- purch. group	rec. quantity	#Orders	#Receipts	#Invoices	Avg order- price	% order- delivery- month=7	% order- delivery- month=8	% order- purch. group=100L50	% order- purch. group=100L51	% order- purch. group=100L52	Avg order- rec. quantity	(Order, % Purchase Other Line Creation)	(Order, % Goods Line Registered)	(Receipt, % Goods Line Registered)	(Invoice, % Invoice Registered)	(Invoice, % Invoice Cleared)
e3	Purchase Requisition Line Created	2017-07-15 12:00																					
e4	Purchase Requisition Line Created	2017-07-15 15:00																					
e5	Purchase Order Line Creation	2017-07-16 15:00																					
e6	Purchase Order Line Creation	2017-07-17 15:00																					
e8	Goods Line Registered	2017-07-22 15:00																					
e9	Invoice Receipt	2017-07-22 16:00																					
e10	Purchase Requisition Group Changed	2017-07-22 19:00																					
e11	Purchase Order Line Creation	2017-07-23 9:00																					
e12	Goods Line Registered	2017-07-23 15:00																					
e13	Invoice Registered	2017-07-29 11:00																					
e14	Invoice Cleared	2017-07-30 12:00																					
e15	Goods Line Registered	2017-07-31 15:00																					
e16	Invoice Registered	2017-08-10 11:00																					
e17	Invoice Cleared	2017-08-15 14:00																					
e4	Purchase Requisition Line Created	2017-07-15 15:00																					
e7	Purchase Order Line Creation	2017-07-18 15:00																					
e18	Goods Line Registered	2017-08-16 15:00																					
e19	Purchase Requisition Supplier Changed	2017-08-16 17:00																					
e20	Invoice Registered	2017-08-18 11:00																					
e21	Invoice Cleared	2017-08-20 14:00																					

Example 3.2. *The result of introducing the aggregation attributes to the feature set is presented in Table 4:*

- *The first type of aggregation attribute that was introduced is represented by the column `Avg order_price`, which represents the average order price (indicated by the column `order_price`) considering all the orders associated to the selected requisition. It can be seen that event `e5` is associated to an `Avg order_price` of 100, since there is only one order with `order_price` 100; conversely, event `e6` is associated to an `Avg order_price` of 150, since there is one order (`o1`) with `order_price` 100 and one order (`o2`) with `order_price` 200.*
- *The second type of aggregation attribute that was introduced is represented by the column `%order_delivery_month=7`, which indicates the fraction of objects correlated with objects of type `order` where attribute `order_delivery_month` takes on value 7: in particular, it indicates the fraction of the orders that are delivered in July (month 7), compared to the total number of orders. As an example, the event `e5` is associated with a `%order_delivery_month=7` of 1, since there is only one order with delivery in July; conversely, event `e6` is associated with a `%order_delivery_month=7` of 0.5, since there is one order (`o1`) with delivery in July and one order (`o2`) with delivery in August (month 8).*
- *The third type of aggregation attribute introduced is represented by the column `#Orders`, which represents the number of objects of type `order` associated to the selected requisition. In the event `e5` there is only one order (`o1`) associated to the selected requisition, while in the event `e6` there are two orders (`o1` and `o2`) associated to the selected requisition.*
- *Finally, the fourth type of aggregation attribute introduced is represented by the column `Order, %Goods Line Registered`, which indicates the percentage of orders that have performed the activity `Goods Line Registered` at least once; as an example, it has a value of 0.5 for event `e8`, since among the*

two orders (o1 and o2) there is only one order (o1) that has performed Goods Line Registered.

4. Implementation and Experiments

Our approach described in Section 3 have been implemented in Python. Prediction models were built using Catboost (cf. Section 2.2). Experiments were conducted on an event log that records the real executions of an object-centric process of an utility provider company in Italy. In particular, we defined five different KPIs of which to predict the values. Section 4.1 introduces the case study employed for our evaluation, while Section 4.2 compares and evaluates proposed predictive techniques in Object-Centric Processes. In particular, we compare our approaches with the naïve adoption of techniques related to single-id process predictive analytics, where only events directly associated with one object type are considered.

4.1. Domain Description

The evaluation of our approach was performed on the object-oriented process described as working example in Section 2. The process is real and is being executed by a well-known Italian utility provider company, which is also one of the major energy companies in Europe. The company focuses on the production/extraction of electricity and gas and on their distribution in different parts of the world. As mentioned in Section 2, the overall process runs through the intertwining of processes related to five different object types.

Section 3 mentioned that our techniques require defining a viewpoint: for the KPIs relevant for this case study, we mainly considered the *Contract* object type as viewpoint, but we also considered the *Order* and the *Requisition* as alternative possible viewpoints. As discussed, if we consider for example the *Contract* viewpoint, traditional predictive analytics based on single-id event logs would only use the events related to *Contract* (the naïve approach). Vice versa, our techniques would also use the information of the objects related to *Contract* via bridge events, including the interaction itself.

Before applying the naïve and our approach, we performed a preprocessing. In particular, we removed attributes whose values were missing in more than 80% of the cases, or whose values were always the same among all cases. In the pre-processing phase, we used domain knowledge and removed the attributes that were somehow duplicate. For instance, the log contains an attribute that refers to the order plant name, which is unique, and a second related to plant identifier: one of the two can be removed. Finally, the large dimension of the utility provider company is also reflected in the cardinality of some categorical attributes. For instance, the codes of the materials that are shipped all around the world (*order_material_code*) are stored in an attribute that counts up to 4,179 different values. To reduce the cardinality of the attributes with thousands of different values, we used the 80-20 rule, a.k.a. Pareto Principle [16]. Specifically, we kept the most frequent attribute’s values that cover the 80% of cases, labelling the remaining values as “other”.

We obtained a different event log for each of the illustrated techniques (naïve approach, approach without aggregated features, approach with aggregated features), consisting of 12,537, 99,065, and 10,349 cases for the *Contract*, the *Order* and the *Requisition* viewpoints respectively.

In each experiment, two thirds of the traces were used for training, and the rest as test set. Splitting was done in a complete random fashion, assuming no concept drift in the dataset. In training, a hyperparameter optimization was carried out, using 20% of the training data for the optimization (validation set).

In our evaluation we considered several KPIs, which can be grouped in three categories. The first is the *path time*, defined as the elapsed time between a defined source activity *a* and the last occurrence of the selected target activity *b*, while the second (*pay delay*) refers to the average number of unforeseen additional days that are needed to receive the payment of the invoice compared to the expected invoice payment due date. Finally, the third category refers to whether or not a certain activity or a certain condition (e.g. a late payment) is eventually going to occur in the future. The first two KPI categories are defined over a numerical domain while the second one is boolean with **true** indicating the occurrence, and **false** the absence.

4.2. Results for Predictive Analytics in Object-Centric Processes

For the evaluation we focused on three different viewpoints (*Contract*, *Requisition* and *Order*) and we considered the two approaches, namely with and without aggregation attributes, which are described in Section 3. Examples 3.1 and 3.2 illustrate the application of our approaches to this case study, when the chosen viewpoint is the *requisition*.

The KPIs that have been considered are reported in Table 5; in particular, Table 5a illustrates the numerical KPIs, while Table 5b illustrates the categorical KPIs. Each numerical KPI refers to the elapsed time from the first occurrence of the considered object to the last occurrence of a different selected activity. The first five KPIs build on the contract viewpoint; in particular, the first KPI refers to the elapsed time from the creation of the contract to the last *SES Line Registered*; this activity indicates that the service requested by the customer has been provided but, since the customer can require several services, it is of interest to know when all the services request have been provided. The second KPI (*SES Line Released*) indicates that a further step has been performed, which is the confirmation from the manager that everything was received correctly. Another interesting KPI to be monitored is the elapsed time from the creation of the contract to the last *Invoice Receipt*, activity that indicates that the invoice has been correctly charged to the customer; conversely, *Invoice Cleared* indicates that the invoice has been paid. The fifth KPI refers to the number of days exceeding the planned payment date (*Pay Delay estimation*), considered starting from the creation of the contract to the last occurrence of *Invoice Cleared*. Finally, the remaining numerical KPIs refer to the elapsed time from the creation of the order or requisition to the last occurrence of some selected activity.

Moreover, after selecting the path from the creation of the contract to the last occurrence of *Invoice Cleared*, we also considered three categorical KPIs, which were related to activities that the company wants to be prevented; in particular, it was interesting to know in advance whether there would be changes to the payment method (represented by the activity *Invoice Pay Method Changed*): when this happens, usually there are delays with payments. Furthermore, the company was also interested in forecasting whether there will be problems with the order (represented by the activity

Table 5: Statistics related to the selected KPIs

(a) Numerical KPIs statistics

KPI	Viewpoint	Average value (days)	Standard deviation value (days)
Elapsed Time from Contract to the last SES Line Registered	Contract	291.59	224.17
Elapsed Time from Contract to the last SES Line Released	Contract	292.73	224.2
Elapsed Time from Contract to the last Invoice Receipt	Contract	257.04	227.41
Elapsed Time from Contract to the last Invoice Cleared	Contract	312.3	240.02
Pay Delay estimation from Contract to the last Invoice Cleared	Contract	9.71	43.2
Elapsed Time from Order to the last Invoice Receipt	Order	24.66	28.08
Elapsed Time from Order to the last Invoice Cleared	Order	64.67	43.45
Elapsed Time from Requisition to the last Invoice Receipt	Requisition	48.17	42.79
Elapsed Time from Requisition to the last Invoice Cleared	Requisition	114	48.56
Elapsed Time from Requisition to the last SES Line Released	Requisition	35.71	40.72
Elapsed Time from Requisition to the last SES Line Registered	Requisition	34.59	40.71

(b) Categorical KPIs statistics. Column % cases represents the percentage of cases in which the activity or the attribute was present with that particular value. Ideally it should be the lowest possible

KPI	Viewpoint	% cases
Occurrence of Activity Pay Method Changed (from Contract to the last Invoice Cleared)	Contract	29%
Occurrence of Activity Purchase Order Blocked (from Contract to the last Invoice Cleared)	Contract	27%
Occurrence of Attribute Pay Type assuming value Late (from Contract to the last Invoice Cleared)	Contract	52%

Purchase Order Blocked), since this situation can bring additional delays caused by the reworks needed to fix the problem. Finally, it was interesting to know whether there will be delays with the payments (represented by the attribute *Pay Type* assuming value Late).

Since eleven KPIs were numerical and the values were reasonably well balanced, we adopted the Mean Absolute Error (MAE). The last three KPIs considered were instead categorical and related to activities that the company wants to be prevented; therefore, we computed the F1 score. Since the selection of the cases and the overall duration of the case itself vary depending on the selected KPI, also the observed average duration and standard deviation of the cases are different; the statistics for the selected numerical KPIs, shown in Table 5a, are necessary to understand if our predictive model is achieving a good prediction quality. Conversely, for the last three categorical KPIs,

Table 6: Prediction performances for the considered KPIs using the proposed techniques, measured in terms of Mean Absolute Error (MAE) or F1 score. The average predictions performances were obtained considering two different divisions of the cases in train and test sets. Training times are reported in brackets.

Viewpoint	KPI	Naïve approach	approach w.out aggr. features	approach w. aggr. features
Contract	Elapsed Time from Contract to the last Invoice Receipt (MAE)	40.67 (5m)	31.72 (36m)	28.49 (53m)
Contract	Elapsed Time from Contract to the last SES Line Released (MAE)	38.63 (7m)	29.96 (25m)	26.71 (15m)
Contract	Elapsed Time from Contract to the last Invoice Cleared (MAE)	44.73 (5m)	37 (33m)	33.06 (81m)
Contract	Elapsed Time from Contract to the last SES Line Registered (MAE)	39.65 (5m)	30.03 (24m)	28.81 (28m)
Contract	Activity Invoice Pay Method Changed occurrence (F1)	0.83 (12m)	0.86 (78m)	0.87 (95m)
Contract	Activity Purchase Order Blocked occurrence (F1)	0.69 (12m)	0.71 (76m)	0.74 (94m)
Contract	Attribute Pay Type Late occurrence (F1)	0.88 (12m)	0.88 (78m)	0.89 (95m)
Contract	Pay Delay estimation from Contract to the last Invoice Cleared (MAE)	15.01 (5m)	14.09 (37m)	13.6 (29m)
Order	Elapsed Time from Order to the last Invoice Receipt (MAE)	11.62 (24m)	9.50 (3h 30m)	10.39 (3h)
Order	Elapsed Time from Order to the last Invoice Cleared (MAE)	15.62 (19m)	12.03 (11h)	11.51 (10h)
Requisition	Elapsed Time from Requisition to the last Invoice Receipt (MAE)	17.65 (7m)	11.43 (31m)	11.7 (33m)
Requisition	Elapsed Time from Requisition to the last Invoice Cleared (MAE)	21.09 (6m)	14.53 (38m)	12.33 (44m)
Requisition	Elapsed Time from Requisition to the last SES Line Released (MAE)	17.56 (5m)	12.84 (11m)	12.79 (12m)
Requisition	Elapsed Time from Requisition to the last SES Line Registered (MAE)	17.35 (5m)	11.3 (14m)	13.08 (14m)

we reported in Table 5b on the distribution of the classes; in particular, the activity *Invoice Pay Method Changed* was performed in the 29% of the cases, the activity *Purchase Order Blocked* was performed in the 27% of cases, while the attribute *Pay Type* appeared with value *Late* in the 52% of cases.

Results of proposed predictive techniques on the selected KPIs are shown in Table 6; in particular, they represent the average accuracy obtained considering two different random splits of the cases in train and test sets. Compared to the naïve approach, it can be seen that the two approaches that consider the object-interaction show improved predictive performances in all the considered KPIs. The additional use of aggregated attributes feeds in additional information that enables a further improvement of the prediction quality in a vast majority of cases.

Table 6 also reports on the training time needed to train the prediction model of the naïve and our approaches on each KPI of interest. As it can be seen, our approaches have the highest training time in all KPIs. Note that this does not pose significant limitations, since it is just performed once; after the model has been trained, all the predictions on the test set can be obtained in less than one second for all the considered KPIs.

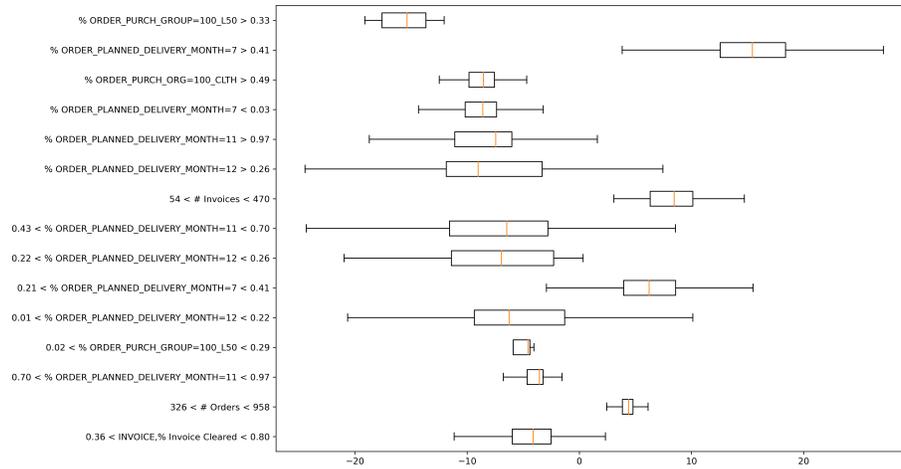


Figure 2: Boxplot representing the impact of some of the most important aggregated features that were added to the predictive model in order to improve the accuracy for the KPI *Elapsed Time from Contract to the last Invoice Receipt* in the approach with aggregated features.

As a further confirmation of the importance of object-interaction and aggregated features, we computed the Shapley Values in accordance to the framework discussed in Section 2.3; in particular, explanations were calculated on the test dataset. Figure 2 illustrates the boxplot representing the distribution of the Shapley values over the different trace prefixes for some of the most important features influencing the prediction of the KPI *Elapsed Time from Contract to the last Invoice Receipt* for the approach with the aggregated features. In particular, each boxplot is ordered by the average Shapley value on the selected KPI considering the absolute value. Each row of the boxplot is linked to an explanation, which extends towards left or right, depending whether the observed Shapley values for the explanation were negative or positive. It can be clearly seen in Figure 2 that many aggregated features are indeed considered relevant by the predictive model; as an example, the most important explanation is *% ORDER_PURCH_GROUP=100_L50 > 0.33* and the average associated Shapley value is -15 days: this means that, when more than the 33% of the orders are related to the purchase group *100_L50*, the estimated Elapsed Time from the signature of the Contract to the last Invoice Receipt reduces on average by 15 days.

A similar reasoning can be applied to Figure 3, which reports on some of the most

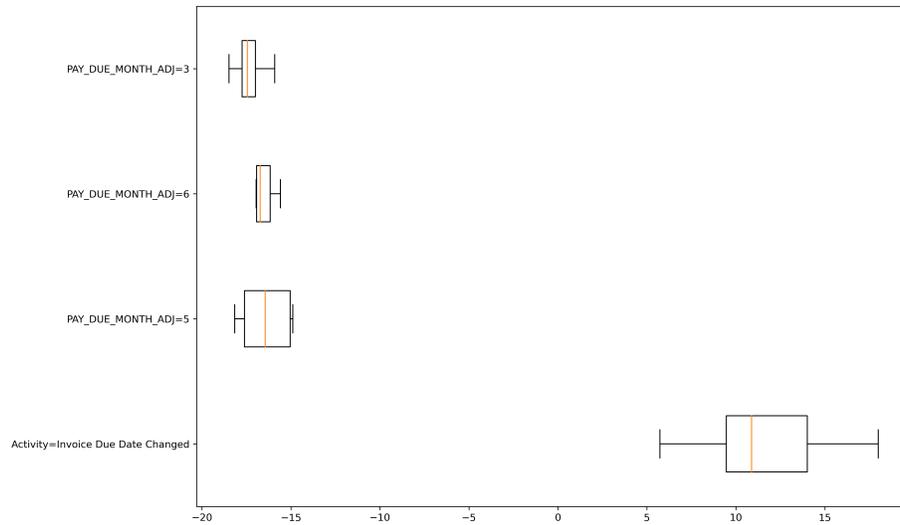


Figure 3: Boxplot representing the impact of some of the most important features that were added to the predictive model (such as the attributes related to the *Invoice* object type) in order to improve the accuracy for the KPI *Elapsed Time from Contract to the last Invoice Cleared* in the approach without aggregated features.

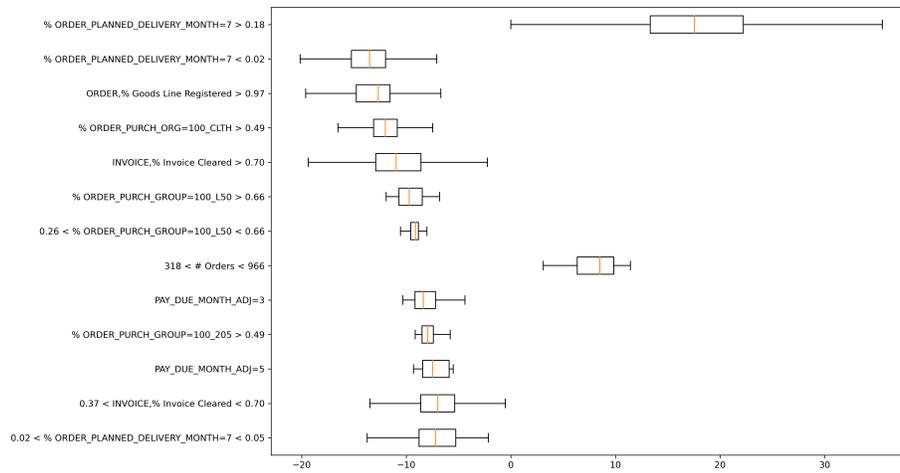


Figure 4: Boxplot representing the impact of some of the most important aggregated features that were added to the predictive model in order to improve the accuracy for the KPI *Elapsed Time from Contract to the last Invoice Cleared* in the approach with aggregated features.

important features influencing the prediction of the KPI *Elapsed Time from Contract to the last Invoice Cleared* for the approach without the aggregated features. As it can be seen, several features were added after considering the object interaction and were considered useful by the predictive model, such as *PAY_DUE_MONTH_ADJ* and the occurrence of the activity *Invoice Due Date Changed*; these are attributes related to the *Invoice* object type, whose values are related to the last opened invoice (or to the last performed activity).

However, since we do not use aggregated features, we miss several factors significantly influencing the prediction. Indeed, when we consider the aggregated features for the same KPI (*Elapsed Time from Contract to the last Invoice Cleared*) and compute the significance of the influence (i.e. the average Shapley value), we obtain the boxplots in Figure 4. It can be clearly seen that the aggregation attributes are now among the most important factors significantly influencing the prediction. As an example, the two most important factors that are contributing to increase the estimated time are represented by $\%Order_planned_delivery_month=7 > 0.18$ and $318 < \#Orders < 966$. The former is associated with an average Shapley value of 20, which means that when more than 18% of the orders are delivered in July (month 7), the estimated time to clear the last invoice is 20 days larger than the average. The latter, associated with an average Shapley value of 8, indicates that when there are a lot of orders that need to be managed at the same time (between 318 and 966), the estimated time to clear the last invoice is 8 days larger than the average. Conversely, one of the most important factors contributing to decrease the estimated time is *ORDER, %Goods Line Registered > 0.97*, which is associated with an average Shapley value of -10; this means that, when the activity Goods Line Registered (which represents the fact that the goods have been received) has been performed at least once in more than 97% of the orders, the estimated time reduces on average by 10 days.

5. Related Works

A body of research exists on object-centric processes. Several research works focus on modelling object-centric processes (e.g. [17]) and the verification of the correctness

of these models [18, 19]. In the realm of Process Mining, techniques are proposed to discover object-centric process models and behavioral dependencies between objects (i.e. artifacts) [20, 21, 22, 23], and to tackle the problem of the object-centric process conformance checking [24, 25]. However, none of the existing works consider object-centric process predictive analytics. A few works consider interactions among different instances of a process [26, 27, 28], but they still rely on the notion of a single process flow (i.e., single case identifier). While some of these works provide valuable insights into inter-case features, their extension to object-centric process is in fact the goal of the technique proposed in this framework.

The work proposed in this paper is based on the idea to unfold an object-centric event log into traditional events logs around a viewpoint. Berti et al. [29] have introduced a very similar concept of viewpoint to indeed extract event logs from the data stored in relational databases: however, the concept of viewpoint is not aimed at process predictive analytics. In fact, unfolding typically leads to a duplication of events, possibly strengthening existing directly-follows relationships or, even, adding new relationships that do not exist in reality. These problems are known as convergence and divergence [2, 30] and make it impossible to apply process-mining techniques designed for single-flow processes that heavily rely on the directly-follow relationships, such as those for model discovery and conformance checking. Conversely, unfolding causes no problem in our process-prediction approach, because we do not use the direct-follow relationships.

6. Conclusions

The lion's share of attention in Business Process Management (BPM) has traditionally been on designing and analyzing processes that are based on a unique notion of case identifier, with a single flow of execution from an initial state to one of the potential final states. In practice, organizations execute more complex processes that are often interacting with each other: one instance of a given process synchronizes with instances of other processes, possibly exchanging data. In light of above, the object-centric process paradigm is nowadays attracting more and more attention in academic

and industry. In this paradigm, the process is seen as the interplay of numerous sub-processes that constitute life cycles of different objects of various types, where these life cycles period synchronize with each other.

This paper tackles the problem of predictive analytics over object-centric processes. The large share of research in predictive analytics cannot be directly applied here, because it traditionally focuses on the problem of predicting the outcome of cases (i.e., process instances) that run in isolation. Also recent techniques that capture the inter-case dynamics assume instances to be of the same process, namely referring to the same object type. If the valuable information of the interaction between objects of the same or different type is not fed into the construction of prediction models, the resulting model might be of low accuracy.

This paper reports on an approach where object-centric logs are unfolded into traditional event logs around a viewpoint, namely around a process' object type. This enables leveraging on the-state-of-art techniques for process predictions.

We conducted experiments with an event log related to a real object-centric process being executed by an utility company in Italy, and measured the accuracy of the predictions using our approach. The accuracy was subsequently also compared with that of a naïve approach that only considers the events related to the process of the viewpoint objects. Experiments have shown that the naïve approach performs rather poorly, confirming the importance of our approach to consider the object interactions when predicting.

As future work, we plan to test on other datasets. We also aim to leverage on prediction techniques based on Graph Neural Networks [31, 32], where the relationships between object types can be explicitly represented as, e.g., a UML or ER diagram (cf. Figure 1). This would likely provide further information for higher-quality predictions.

References

- [1] G. Li, R. Medeiros de Carvalho, W. M. P. van der Aalst, Object-centric behavioral constraint models: A hybrid model for behavioral and data perspectives, in:

Proceedings of the 34th ACM/SIGAPP Symposium on Applied Computing (SAC 2019), ACM, 2019, p. 48–56.

- [2] W. M. P. van der Aalst, Object-centric process mining: Dealing with divergence and convergence in event data, in: Proceedings of 17th International Conference on Software Engineering and Formal Methods (SEFM 2019), Vol. 11724 of LNCS, Springer, 2019, pp. 3–25.
- [3] W. M. P. van der Aalst, A. Artale, M. Montali, S. Tritini, Object-centric behavioral constraints: Integrating data and declarative process modelling, in: Description Logics, Vol. 1879 of CEUR Workshop Proceedings, CEUR-WS.org, 2017.
- [4] M. T. Wynn, J. Leberherz, W. M. P. van der Aalst, R. Accorsi, C. Di Ciccio, L. Jayarathna, H. Verbeek, Rethinking the input for process mining: Insights from the xes survey and workshop, in: Proceedings of the ICPM 2021 Process Mining Workshops, LNBIP, Springer, in press. Preprint available at <https://www.tf-pm.org/upload/1637654003748.pdf>.
- [5] A. Farhang Ghahfarokhi, G. Park, A. Berti, W. M. P. van der Aalst, OCEL Standard, Available at <http://ocel-standard.org/1.0/specification.pdf> (January 2020).
- [6] A. E. Márquez-Chamorro, M. Resinas, A. Ruiz-Cortés, Predictive monitoring of business processes: A survey, IEEE Transaction on Services Computing 11 (6) (2018) 962–977.
- [7] G. Park, M. Song, Prediction-based resource allocation using LSTM and minimum cost and maximum flow algorithm, in: Proceedings of the International Conference on Process Mining, ICPM 2019, Aachen, Germany, 2019, IEEE, 2019, pp. 121–128.
- [8] N. Tax, I. Verenich, M. La Rosa, M. Dumas, Predictive business process monitoring with LSTM neural networks, in: Proceedings of the 29th International Conference on Advanced Information Systems Engineering (CAiSE 2017), Vol. 10253, Springer, 2017, pp. 477–492.

- [9] N. Navarin, B. Vincenzi, M. Polato, A. Sperduti, LSTM networks for data-aware remaining time prediction of business process instances, in: Proceedings of the IEEE Symposium Series on Computational Intelligence (SSCI 2017), 2017, pp. 1–7.
- [10] M. Polato, A. Sperduti, A. Burattin, M. de Leoni, Time and activity sequence prediction of business process instances, *Computing* 100 (9) (2018) 1005–1031.
- [11] I. Verenich, M. Dumas, M. La Rosa, F. Maggi, I. Teinmaa, Survey and cross-benchmark comparison of remaining time prediction methods in business process monitoring, *ACM Transactions on Intelligent Systems and Technology* 10 (2019) 1–34.
- [12] A. V. Dorogush, V. Ershov, A. Gulin, Catboost: gradient boosting with categorical features support, in: Proceedings of the Workshop on ML Systems at NIPS 2017, 2017.
- [13] R. Galanti, B. Coma-Puig, M. de Leoni, J. Carmona, N. Navarin, Explainable predictive process monitoring, in: Proceedings of the 2nd International Conference on Process Mining (ICPM 2020), IEEE, 2020, pp. 1–8.
- [14] L. S. Shapley, A value for n-person games, Vol. 2, RAND Corporation, 1953, pp. 307–317.
- [15] C. Molnar, G. Casalicchio, B. Bischl, Interpretable machine learning - A brief history, state-of-the-art and challenges, in: Proceedings of workshops of the European Conference on Machine Learning and Knowledge Discovery in Databases (ECML PKDD 2020), Ghent, Belgium, September 14-18, 2020, Vol. 1323 of Communications in Computer and Information Science, Springer, 2020, pp. 417–431.
- [16] M. E. J. Newman, Power laws, Pareto distributions and Zipf’s law, Vol. 46, Taylor & Francis, 2005, pp. 323–351.

- [17] D. Cohn, R. Hull, Business artifacts: A data-centric approach to modeling business operations and processes, *IEEE Data Engineering Bulletin* 32 (2) (2009) 3–9.
- [18] D. Calvanese, M. Montali, M. Estañol, E. Teniente, Verifiable uml artifact-centric business process models, in: *Proceedings of the 23rd ACM International Conference on Conference on Information and Knowledge Management, CIKM '14*, Association for Computing Machinery, New York, NY, USA, 2014, p. 1289–1298.
- [19] M. Montali, D. Calvanese, Soundness of data-aware, case-centric processes, *International Journal on Software Tools for Technology Transfer* 18 (5) (2016) 535–558.
- [20] W. M. P. van der Aalst, A. Berti, Discovering object-centric petri nets, *Fundamenta Informaticae* 175 (1-4) (2020) 1–40.
- [21] X. Lu, M. Nagelkerke, D. van de Wiel, D. Fahland, Discovering Interacting Artifacts from ERP Systems, *IEEE Transactions on Services Computing* 8 (6) (2015) 861–873. doi:10.1109/TSC.2015.2474358.
- [22] V. Popova, M. Dumas, Discovering unbounded synchronization conditions in artifact-centric process models, in: *Proceedings of Business Process Management Workshops*, Vol. 171 of LNBIP, Springer, 2014, pp. 28–40.
- [23] M. L. van Eck, N. Sidorova, W. M. P. van der Aalst, Guided interaction exploration in artifact-centric process models, in: *Proceedings of the 19th Conference on Business Informatics (CBI)*, Vol. 01, IEEE, 2017, pp. 109–118. doi:10.1109/CBI.2017.42.
- [24] D. Fahland, M. de Leoni, B. F. van Dongen, W. M. P. van der Aalst, Conformance checking of interacting processes with overlapping instances, in: *Proceedings of the 9th International Conference on Business Process Management (BPM 2011)*, Vol. 6896, Springer Berlin Heidelberg, 2011, pp. 345–361.

- [25] J. N. Adams, W. M. P. van der Aalst, Precision and fitness in object-centric process mining, in: *Proceedings of the 3rd International Conference on Process Mining (ICPM 2021)*, IEEE, 2021, pp. 128–135.
- [26] A. Senderovich, C. Di Francescomarino, F. Maggi, From knowledge-driven to data-driven inter-case feature encoding in predictive process monitoring, *Inf. Syst.* 84 (2019) 255–264.
- [27] V. Denisov, D. Fahland, W. M. P. van der Aalst, Predictive performance monitoring of material handling systems using the performance spectrum, in: *Proceedings of the International Conference on Process Mining, ICPM 2019*, 2019, pp. 137–144. doi:10.1109/ICPM.2019.00029.
- [28] E. L. Klijn, D. Fahland, Identifying and reducing errors in remaining time prediction due to inter-case dynamics, in: *Proceedings of the 2nd International Conference on Process Mining, ICPM 2020*, 2020, pp. 25–32. doi:10.1109/ICPM49681.2020.00015.
- [29] W. M. P. Alessandro Berti and, van der Aalst, Extracting multiple viewpoint models from relational databases, in: *Proceedings of the 9th International Symposium on Data-driven Process Discovery and Analysis (SIMPDA 2019)*, Vol. 379 of *Lecture Notes in Business Information Processing*, Springer, 2019, pp. 24–51.
- [30] S. Esser, D. Fahland, Multi-dimensional event data in graph databases, *Journal on Data Semantics* 10 (1) (2021) 109–141. doi:10.1007/s13740-021-00122-1.
- [31] F. Scarselli, M. Gori, A. C. Tsoi, M. Hagenbuchner, G. Monfardini, The graph neural network model, *IEEE Transactions on Neural Networks* 20 (1) (2009) 61–80.
- [32] F. Wu, A. Souza, T. Zhang, C. Fifty, T. Yu, K. Weinberger, Simplifying graph convolutional networks, in: *Proceedings of the 36th International Conference on Machine Learning (ICML) 2019*, Vol. 97, PMLR, 2019, pp. 6861–6871.