

Feature Recommendation for Structural Equation Model Discovery in Process Mining

Mahnaz Sadat Qafari and Wil van der Aalst

Rheinisch-Westfälische Technische Hochschule Aachen(RWTH), Aachen, Germany
m.s.qafari@pads.rwth-aachen.de, wvdaalst@pads.rwth-aachen.de

Abstract. Process mining techniques can help organizations to improve the operational processes. Organizations can benefit from process mining techniques in finding and amending the root causes of performance or compliance problems. Considering the volume of the data and the number of features captured by the information system of today's companies, the task of discovering the set of features that should be considered in root cause analysis can be quite involving. In this paper, we propose a method for finding the set of (aggregated) features with a possible effect on the problem. The root cause analysis task is usually done by applying a machine learning technique to the data gathered from the information system supporting the processes. To prevent mixing up correlation and causation, which may happen because of interpreting the findings of machine learning techniques as causal, we propose a method for discovering the structural equation model of the process that can be used for root cause analysis. We have implemented the proposed method as a plugin in ProM and we have evaluated it using two real and synthetic event logs. These experiments show the validity and effectiveness of the proposed methods.

Keywords: Process mining · Root cause analysis · Causality inference.

1 Introduction

Organizations aim to improve operational processes to serve customers better and to become more profitable. To this goal, they can utilize process mining techniques in many steps, including identifying friction points in the process, finding the root causes of each friction point, estimating the possible impact of changing each factor on the process performance, and also planning process enhancement actions. Today, there are several robust techniques for process monitoring and finding their friction points, but little work on *root cause analysis*. So, in this paper, we focus on root cause analysis and investigating the impact of interventions.

Processes are complicated entities involving many steps, where each step itself may include many influential factors and features. Moreover, not just the steps but also the order of the steps that are taken for each process instance may vary, which results in several process instance variants. This makes it quite hard to identify the set of features that influence a problem. In the literature related to root cause analysis in the field of process mining, it is usually assumed that the user provides the set of features that have a causal relationship with the observed problem in the process (see for example [9,16]).

To overcome this issue, we have proposed a mechanism that not only helps to identify the set of features that may have a causal relationship with the problem, but also the values of these features that are more prone to causing the problem.

Traditionally, the task of finding the root cause of a problem in a process is done in two steps; first gathering process data from the event log, and then applying data mining and machine learning techniques. Although the goal is to perform root cause analysis, a naive application of such techniques often leads to a mix-up of correlation and causation. It is easy to find correlations, but very hard to determine causation. Consequently, process enhancement based on the results of such approaches does not always lead to any process improvements.

Consider the following three scenarios:

- (i) In an online shop, it has been observed that the possibility of delay in delivery is much higher if some specific resources are responsible for them.
- (ii) In a consultancy company, there are deviations in some cases and those cases have been done mainly by the employees who are most experienced.
- (iii) In an IT company, it has been observed that the higher the number of resources assigned to a task, the longer it takes.

The following possibly incorrect conclusions can be made if these observed correlations are observed as causal relationships.

- In the online shop scenario, the responsible resources are causing the delays.
- In the second scenario, we may conclude that over time the employees get more and more reckless, and consequently the rate of deviations increases.
- In the IT company, we may conclude that the more people working on a project, the more time is spent on team management and communication, which prolongs the project unnecessarily.

However, correlation does not mean causation. We can have a high correlation between two events when they are caused by a possibly unmeasured (hidden) common cause (set of common causes), which is called a *confounder*. For example, in the first scenario, the delayed deliveries are mainly for the bigger size packages which are usually assigned to specific resources. Or, in the second scenario, the deviations happen in the most complicated cases that are usually handled by the most experienced employees. In the third scenario, maybe both the number of employees working on a project and the duration of a project are highly dependent on the complexity of the project. As it is obvious from these examples, changing the process based on the observed correlations not only leads to no improvement but also may aggravate the problem (or create new problems).

Randomized experiments and the theory of causality are two general frameworks for finding the causes of a problem [13,14]. The randomized experiment provides the most robust and reliable method for making causal inferences and statistical estimates of the effect of an intervention. This method involves randomly setting the values of the features that have a causal effect on the observed problem and monitoring the effects. Applying randomized experiments in the context of processes is usually too expensive (and sometimes unethical) or simply impossible. The other option for anticipating the effect of any intervention on the process is using a *structural causal model* [13,14].

In this method, first, the causal mechanism of the process features is modeled by a conceptual model and then this model is used for studying the effect of changing the value of a process features.

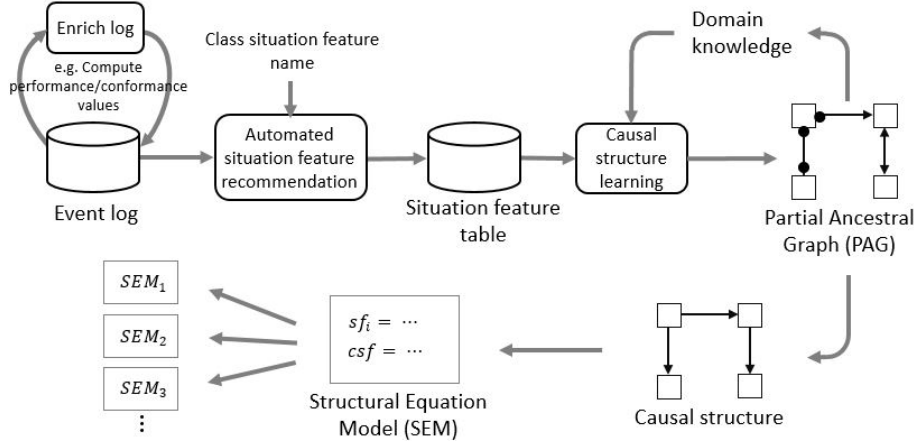


Fig. 1: The general structural causal equation discovery.

This paper is an extension of [15], where we have proposed a framework for root cause analysis using structural equation modeling. Here we address one of the main issues in this framework. Finding the features that may have a causal effect on the problem often requires substantial domain knowledge. Considering the variety of the feature values in a process, even when having extensive domain knowledge, it may not be easy to determine the values of the features that have the strongest effect on the problem. So, we propose a method for finding a set of features and feature value pairs that may contribute the most to the problem. Also, we add aggregated features to the features that can be extracted from the event log. This makes the method capable of analyzing more scenarios. The framework explained in this paper includes the following steps:

- As a preprocessing step, the event log is enriched by several process-related features. These features are derived from different data sources like the event log, the process model, and the conformance checking results. Also, here we consider the possibility of adding aggregated features to the event log regarding the time window provided by the user.
- A set of pairs of the form (feature, feature value) are recommended to the user. Such pair include the features that might have a causal relationship with the problem and those values of them that possibly contribute more to the problem. Users can modify this set of features that have been identified automatically or simply ignore it and provide another set of features.
- The next step is creating a target-dependent data table, which we call it *situation feature table*.

- This step involves generating a graphical object encoding the structure of causal relationships among the process features. This graphical object can be provided by the customer or be inferred from the observational data using a causal structure learning algorithm, also called *search algorithm*. The user can modify the resulting graphical object by adding domain knowledge as an input to the search algorithm or by modifying the discovered graph.
- The last step involves estimating the strength of each discovered causal relationship and the effect of an intervention on any of the process features on the identified problem.

In Figure 1, the general overview of the proposed approach is presented.

The remainder of the paper is organized as follows. In Section 2, we start with an example. We use this example as the running example throughout this paper. In Section 3, we present some of the related work. The corresponding process mining and causal inference theory preliminaries are presented in Section 4 and, in Section 5, an overview of the proposed approaches for feature recommendation and causal equation model discovery is presented. In Section 6, the assumptions and the design choices in the implemented plugin and the experimental results of applying it on synthetic and real event logs are presented. Finally, in Section 7, we summarize our approach and its applications.

2 Motivating Example

As the running example, we use an imaginary IT company that implements software for its customers. However, they do not do the maintenance of the released software. Here, each process instance is corresponding to the process of implementing one software. This process involves the following activities: business case development, feasibility study, product backlog, team charter, development, test, and release. The Petri-net model of this company is shown in Figure 2. We refer to the sub-model including two transitions “development” and “test” (the two blue activities in Figure 2) as *implementation phase*.

The manager of the company is concerned about the duration of the implementation phase of projects. She wants to know what features determine the implementation phase duration. And also, if there is any way to reduce the implementation phase duration. If so, what would be the effect of changing each feature. These are valid questions to be asked before planning for re-engineering and enhancing the process. The manager believes that the following features of a project are the process features that might have a causal effect on its “*implementation phase duration*” (the duration of implementation phase):

- “*Priority*” which is an attribute of business case development indicating how urgent the software is for the customer,
- “*Team size*” which is an attribute of team charter indicating the number of resources working on a project,
- “*Duration*” of product backlog activity, an attribute of product backlog, which indicates the duration of the product backlog activity.

Analyzing the historical data from the company shows that there is a high correlation between every one of the three mentioned features and the duration of the implementation phase. We consider “Complexity” (the complexity and hardness of the project) as another feature that is not recorded in the event log but has a causal effect on the duration of the implementation phase.

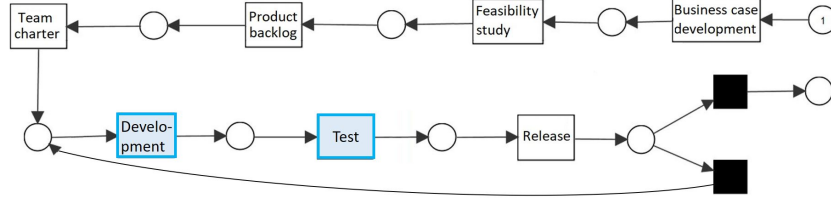


Fig. 2: The Petri net model of the process of IT company described in Section 2.

The structure of the causal relationship among the features has a high impact on the answers to the mentioned questions. In Figures 3, 4, and 5, three possible structures of the causal relationship among the features of the IT company are depicted¹.

According to Figure 3, just team size and priority have a causal effect on the duration of the implementation phase. But product backlog duration does not have any causal effect on the duration of the implementation phase even though they are highly correlated. Consequently, changing product backlog duration does not have any impact on the duration of the implementation phase.

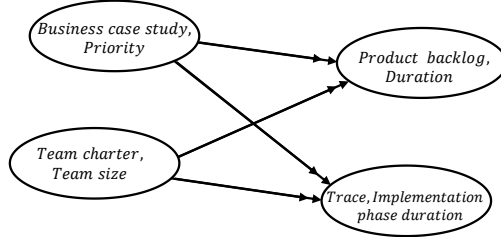


Fig. 3: A possible causal structure for the IT company.

According to Figure 4, all three features priority, product backlog duration, and team size influence the duration of the implementation phase. Thus, by changing each of these three features, one can influence the duration of the implementation phase.

¹ In these three figures and other figures in this paper that visualize networks of feature, the labels of the nodes are either of the form *Trace, Attribute name* if the attribute name is related to a trace-level attribute, or of the form *Activity name, Attribute name* if the attribute name is related to an event-level attribute. In the former case, the activity name indicated the activity that the attribute belongs to.

Based on 5, we can conclude that the complexity, which is a hidden feature in the model (depicted by the gray dashed oval in Figure 5), causally influences both implementation phase duration and product backlog duration. Subsequently, the correlation among them is because of having a common cause. Grounded in this causal structure, it is not possible to influence the duration of the implementation phase by forcing product backlog activity to take place in a shorter or longer amount of time.

It is worth noting that not all the features are actionable, i.e., in reality, it is not possible to intervene on some of the features. For example, in the mentioned IT company, we can imagine that the manager intervenes on team size by assigning more or fewer people to a project; but he cannot intervene in the complexity of a project. Judging whether a feature can be intervened requires using common sense and domain knowledge.

In the rest of this paper, we show how to answer such questions posed by the company manager. We first mention how to extract data in a meaningful way regarding the target feature (implementation phase duration in this example) and then we show how to discover the causal relationships between the process features and the structural equation model of the features that may affect the target feature using our method. Finally, we demonstrate how questions related to investigating the effect of intervention on the target feature can be answered in this framework. In Section 6.2, we show the results of applying our method for answering the mentioned question by the IT company manager in this example.

3 Related Work

In the literature, there is plenty of work in the area of process mining dedicated to finding the root causes of a performance or compliance problem. The root cause analysis approach of the proposed methods usually involves classification [4,9], and rule mining

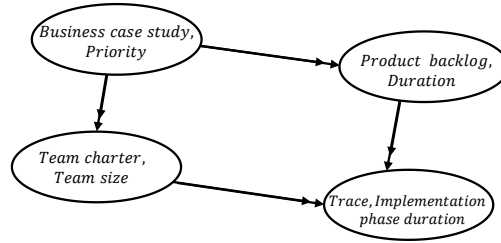


Fig. 4: A possible causal structure for the IT company.

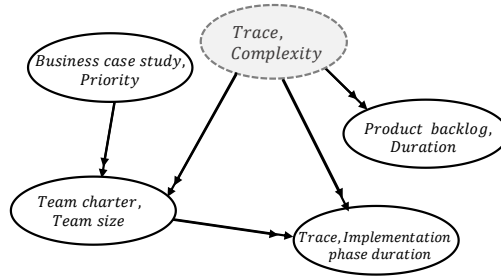


Fig. 5: A possible causal structure for the IT company.

[18]. The main problem of these approaches is that the findings of these methods are based on correlation which does not necessarily imply causation.

The theory of causation based on the structural causal model has been studied deeply [14]. Also, a variety of domains benefit from applying methods from this domain (e.g. [10,22]). However, there is little work on the application of the theory of causality in the area of process mining. There are some works in process mining that use causality theory. These includes:

- In [5], the authors propose an approach for discovering causal relationships between a range of business process characteristics and process performance indicators based on time-series analysis. The idea is to generate a set of time-series using the values of performance indicators, and then applying Granger causality test on them, to investigate and discover their causal relationships. Granger test is a statistical hypothesis test to detect predictive causality; consequently, the causal relationships using this approach might not be true cause-and-effect relationships.
- In [11], the authors use the event log and the BPMN model of a process to discover the structural causal model discovery of the features of the process. They first apply loop unfolding on the BPMN model of the process and generate a partial order of features. They use the generated partial order to guide the search algorithm. In this work, it is assumed that the BPMN model of a process is its accurate model, which is not always the case.

There is also some work devoted to the case level root cause analysis [1,16].

It is worth mentioning that all the above-mentioned approaches are based on statistical tests for discovering causal relationships. Consequently, these approaches are not feasible when there are a huge number of features. However, none of them provides a method for feature recommendation. Yet, there is some work on influence analysis that aims at finding such a set of features [6,7,8].

4 Preliminaries

In this section, we describe some of the basic notations and concepts of the process mining and causal inference theory.

In the following section, we follow two goals: first, we describe the basic notations and concepts of the process mining and second, we show the steps involved in converting a given event log into a situation feature table.

4.1 Process Mining

Process mining techniques start from an *event log* extracted from an information system. The atomic building block of an event log is an *event*. An event indicates that an activity happened at a specific point in time for a specific case. A set of events that are related to a specific process instance are called a *trace*. We can look at an event log as a collection of traces. An event log may include three different levels of attributes: log-level attributes, trace-level attributes, and event-level attributes. In the following, we explicitly define an event, trace and event log in a way that reflects reality and at

the same time is suitable for our purpose. But first, we need to define the following universes and functions:

- \mathcal{U}_{att} is the universe of *attribute names*, where $\{actName, timestamp, caseID\} \subseteq \mathcal{U}_{att}$. *actName* indicates the activity name, *timestamp* indicates the timestamp of an event, and *caseID* indicates that the event belongs to which case.
- \mathcal{U}_{val} is the universe of *values*.
- $values \in \mathcal{U}_{att} \mapsto \mathbb{P}(\mathcal{U}_{val})$ as a function that returns the set of all possible values of a given attribute name².
- $\mathcal{U}_{map} = \{m \in \mathcal{U}_{att} \mapsto \mathcal{U}_{val} \mid \forall at \in dom(m) : m(at) \in values(at)\}$ the universe of all mappings from a set of attribute names to attribute values of the correct type.

Also, we define \perp as a member of \mathcal{U}_{val} such that $\perp \notin values(at)$ for all $at \in \mathcal{U}_{att}$. We use this symbol to indicate that the value of an attribute is unknown, undefined, or is missing.

Now, we define an event as follows:

Definition 1 (Event). *An event is an element of $e \in \mathcal{U}_{map}$, where $e(actName) \neq \perp$, $e(timestamp) \neq \perp$, and $e(caseID) \neq \perp$. We denote the universe of all possible events by \mathcal{E} and the set of all non-empty chronologically ordered sequences of events that belong to the same case (have the same value for *caseID*) by \mathcal{E}^+ . If $\langle e_1, \dots, e_n \rangle \in \mathcal{E}^+$, then for all $1 \leq i < j \leq n$, $e_i(timestamp) \leq e_j(timestamp) \wedge e_i(caseID) = e_j(caseID)$.*

Example 1. The events in the following table are some of the possible events for the IT company in Section 2.

$e_1 := \{(caseID, 1), (actName, \text{“Business case development”}), (timestamp, t_1), (Priority, 2)\}$
$e_2 := \{(caseID, 1), (actName, \text{“Feasibility study”}), (timestamp, t_2)\}$
$e_3 := \{(caseID, 1), (actName, \text{“Product backlog”}), (timestamp, t_3), (Duration, 35)\}$
$e_4 := \{(caseID, 1), (actName, \text{“Team charter”}), (timestamp, t_4), (team size, 21)\}$
$e_5 := \{(caseID, 1), (actName, \text{“Development”}), (timestamp, t_5), (Duration, 200)\}$
$e_6 := \{(caseID, 1), (actName, \text{“Test”}), (timestamp, t_6), (Duration, 79)\}$
$e_7 := \{(caseID, 1), (actName, \text{“Release”}), (timestamp, t_7)\}$
$e_8 := \{(caseID, 2), (actName, \text{“Business case development”}), (timestamp, t_8), (Priority, 1)\}$
$e_9 := \{(caseID, 2), (actName, \text{“Feasibility study”}), (timestamp, t_9)\}$
$e_{10} := \{(caseID, 2), (actName, \text{“Product backlog”}), (timestamp, t_{10}), (Duration, 63)\}$
$e_{11} := \{(caseID, 2), (actName, \text{“Team charter”}), (timestamp, t_{11}), (team size, 33)\}$
$e_{12} := \{(caseID, 2), (actName, \text{“Development”}), (timestamp, t_{12}), (Duration, 226)\}$
$e_{13} := \{(caseID, 2), (actName, \text{“Test”}), (timestamp, t_{13}), (Duration, 74)\}$
$e_{14} := \{(caseID, 2), (actName, \text{“Release”}), (timestamp, t_{14})\}$
$e_{15} := \{(caseID, 2), (actName, \text{“Development”}), (timestamp, t_{15}), (Duration, 62)\}$
$e_{16} := \{(caseID, 2), (actName, \text{“Test”}), (timestamp, t_{16}), (Duration, 117)\}$
$e_{17} := \{(caseID, 2), (actName, \text{“Release”}), (timestamp, t_{17})\}$

² In this paper, it is assumed that the reader is familiar with sets, multi-sets, and functions. $\mathbb{P}(X)$ is the set of non-empty subsets of set $X \neq \emptyset$. Let X and Y be two sets. $f : X \mapsto Y$ is a partial function. The domain of f is a subset of or equal to X which is denoted by $dom(f)$. We write $f(x) = \perp$ if $x \notin dom(f)$.

Each event may have several attributes which can be used to group the events. For $at \in \mathcal{U}_{att}$, and $V \subseteq values(at)$, we define a group of events as the set of those events in \mathcal{E} that assign a value of V to the attribute at ; i.e.

$$group(at, V) = \{e \in \mathcal{E} | e(at) \in V\}.$$

Some of the possible groups of events are:

- the set of events with specific activity names,
- the set of events which are done by specific resources,
- the set of events that start in a specific time interval during the day, or,
- the set of events with a specific duration.

We denote the universe of all event groups by $\mathcal{G} = \mathbb{P}(\mathcal{E})$.

Example 2. Here are some possible event groups based on the IT company in Section 2.

$$\begin{aligned} G_1 &:= group(actName, \{\text{“Business case development”}\}) \\ G_2 &:= group(actName, \{\text{“Product backlog”}\}) \\ G_3 &:= group(actName, \{\text{“Team charter”}\}) \\ G_4 &:= group(actName, \{\text{“Development”}\}) \\ G_5 &:= group(team\ size, \{33, 34, 35\}) \end{aligned}$$

Based on the definition of an event, we define an event log as follows:

Definition 2 (Event Log). We define the universe of all event logs as $\mathcal{L} = \mathcal{E}^+ \not\rightarrow \mathcal{U}_{map}$. Let L where $L \in \mathcal{L}$ be an event log, we call each element $(\sigma, m) \in L$ a trace.

One of our assumptions in this paper is the uniqueness of events in event logs; i.e., given an event log $L \in \mathcal{L}$, we have $\forall (\sigma_1, m_1), (\sigma_2, m_2) \in L : e_1 \in \sigma_1 \wedge e_2 \in \sigma_2 \wedge e_1 = e_2 \implies (\sigma_1, m_1) = (\sigma_2, m_2)$ and $\forall (\langle e_1, \dots, e_n \rangle, m) \in L : \forall 1 \leq i < j \leq n : e_i \neq e_j$. This property can easily be ensured by adding an extra identity attribute to the events.

Also, we assume that the uniqueness of the “caseID” value for traces in a given event log L . In other words, $\forall (\sigma_1, m_1), (\sigma_2, m_2) \in L : e_1 \in \sigma_1 \wedge e_2 \in \sigma_2 \wedge e_1(caseID) = e_2(caseID) \implies (\sigma_1, m_1) = (\sigma_2, m_2)$.

Example 3. $L_{IT} = \{\lambda_1, \lambda_2\}$ is a possible event log for the IT company in 2. L_{IT} includes two traces λ_1 and λ_2 , where:

- $\lambda_1 := (\langle e_1, \dots, e_7 \rangle, \{(Responsible, Alice)\})$ and
- $\lambda_2 := (\langle e_8, \dots, e_{17} \rangle, \{(Responsible, Alex)\})$.

Here t_1, \dots, t_{17} are unique timestamps where $t_1 < \dots < t_7$ and $t_8 < \dots < t_{17}$.

As a preprocessing step, we enrich the event log by adding many derived features to its traces and events. There are many different derived features related to any of the process perspectives; the time perspective, the data flow-perspective, the control-flow perspective, the conformance perspective, or the resource/organization perspective of

the process. We can compute the value of the derived features from the event log or possibly other sources.

Moreover, we can enrich the event log by adding aggregated attributes to its events and traces. Let $L \in \mathcal{L}$ be an event log, $k \in \mathbb{N}$ (a non-zero natural number) the number of time windows, t_{min} the minimal timestamp, and t_{max} the maximum timestamp in L , we divide the time span of L into k consecutive time windows with equal length (the length of each time window is $(t_{max} - t_{min})/k$ and compute the value of aggregated attributes for each of these k time windows. In other words, We define $\xi : \mathcal{L} \times \mathcal{U}_{att} \times \mathbb{N} \times values(timestamp) \rightarrow \mathbb{R}$ as a function that given an event log, an aggregated attribute name, the number of time windows, and a timestamp returns the value of the given aggregated attribute in the time window that includes the timestamp. We can use ξ for aggregated attributes at both the event and the trace-level. More precisely, given $L \in \mathcal{L}$, $(\sigma, m) \in L$, $e \in \sigma$, $k \in \mathbb{N}$, and $at \in \mathcal{U}_{att}$ where at is an aggregated attribute, we define $e(at) = \xi(L, at, k, e(timestamp))$ and $m(at) = \xi(L, at, k, t')$ where $t' = \max\{e(timestamp) | e \in \sigma\}$. Some of the possible aggregated attributes are: the number of waiting customers, workload (in the process-level), average service time, average waiting time (in the trace and event-level), number of active events with a specific activity name, number of waiting events with a specific activity name (in the event-level), average service time, average waiting time (in the resource-level).

While extracting the data from an event log, we assume that the event recording delays by the information system of the process were negligible. Considering the time order of cause and effect, we have that only the features that have been recorded before the occurrence of a specific feature can have a causal effect on it. So the relevant part of a trace to a given feature is a prefix that trace, which we call such a prefix of a trace a *situation*. Let $prfx(\langle e_1, \dots, e_n \rangle) = \{\langle e_1, \dots, e_i \rangle | 1 \leq i \leq n\}$, a function that returns the set of non-empty prefixes of a given sequence of events. Using $prfx$ function we define a situation as follows:

Definition 3 (Situation). We define $\mathcal{U}_{situation} = \mathcal{E}^+ \times \mathcal{U}_{map}$ as the universe of all situations. We call each element $(\sigma, m) \in \mathcal{U}_{situation}$ a situation. Considering $L \in \mathcal{L}$, we define the set of situations of L as $S_L = \{(\sigma, m) | \sigma \in prfx(\sigma') \wedge (\sigma', m) \in L\}$.

Among the possible subsets of S_L of a given event log L , we distinguish two important type situation subsets of S_L . The first type is the *G-based situation subset* of L where $G \in \mathcal{G}$ and includes those situations in S_L that their last event (the event with maximum timestamp) belongs to G . The second type is the *trace-based situation subset*, which includes the set of all traces of L .

Definition 4 (Situation Subset). Given $L \in \mathcal{L}$ where $S_L \subseteq \mathcal{U}_{situation}$ is the set of situations of L , and $G \in \mathcal{G}$, we define

- *G-based situation subset* of L as $S_{L,G} = \{(\langle e_1, \dots, e_n \rangle, m) \in S_L | e_n \in G\}$, and
- *trace-based situation subset* of L as $S_{L,\perp} = L$.

Example 4. Three situations s_1 , s_2 , and s_3 , where $s_1, s_2, s_3 \in S_{LIT, G_4}$ (G_4 in Example 2, generated using the trace mentioned in Example 3) are as follows:

$$\begin{aligned} s_1 &:= (\langle e_1, \dots, e_5 \rangle, \{(Responsible, Alice)\}) \\ s_2 &:= (\langle e_8, \dots, e_{12} \rangle, \{(Responsible, Alex)\}) \\ s_3 &:= (\langle e_8, \dots, e_{15} \rangle, \{(Responsible, Alex)\}) \end{aligned}$$

Note that $G_4 := group(actName, \{\text{“Development”}\})$ and we have $\{e_5, e_{12}, e_{15}\} \subseteq G_4$. In other words $e_5(actName) = e_{12}(actName) = e_{15}(actName) = \text{“Development”}$.

If a process includes decision points, then one of the derived attributes that can be added to the event log when enriching the event log is the *choice* attribute. A choice attribute is added to the activity that happens before the decision point and its value indicates which activity has been enabled as the result of the decision that has been made. So we can use an added choice attribute and its values to group the events in an event log and extract a situation subset based on the occurrence of that specific choice. We already defined two important types of situation subsets; group-based situation subsets and trace-based situation subsets. We also distinguish the *choice-based situation subsets* where the situation subset is extracted based on events that have a specific choice attribute. These situation subsets are important as they are conceptually related to a decision point.

When extracting the data, we need to distinguish trace-level attributes from event-level attributes. We do that by using *situation features* which is identified by a group of events, G (possibly $G = \perp$), and an attribute name, at . Each situation feature is associated with a function defined over the situations. This function returns the proper value for the situation feature regarding at and G extracted from the given situation. More formally:

Definition 5 (Situation Feature). We define $\mathcal{U}_{sf} = \mathcal{U}_{att} \times (\mathcal{G} \cup \{\perp\})$ as the universe of the situation features. Each situation feature is associated with a function $\#_{sf} : \mathcal{U}_{situation} \mapsto \mathcal{U}_{val}$ such that given $sf = (at, G)$ where $at \in \mathcal{U}_{att}$, and $G \in \mathcal{G} \cup \{\perp\}$

we have:

- if $G = \perp$, then $\#_{(at, G)}((\sigma, m)) = m(at)$ and
- if $G \in \mathcal{G}$, then $\#_{(at, G)}((\sigma, m)) = e(at)$ where $e = \arg \max_{e' \in G \cap \{e'' \in \sigma\}} e'(timestamp)$ for $(\sigma, m) \in \mathcal{U}_{situation}$.

We denote the universe of the situation features as \mathcal{U}_{sf} .

We can consider a situation feature as an analogy to the feature (a variable) in a tabular data. Also, we can look at the corresponding function of a situation feature as the function that determines the mechanism of extracting the value of the situation feature from a given situation. Given a situation (σ, m) and a situation feature (at, G) , if $G = \perp$, its corresponding function returns the value of at in trace level (i.e., $m(at)$). However, if $G \neq \perp$, then the function returns the value of at in $e \in \sigma$ that belongs to G and happens last (has the maximum timestamp) among those events of σ that belong to G .

Example 5. We can define the following situation features using the information mentioned in the previous examples:

$$\begin{aligned} sf_1 &:= (Team\ size, G_3) & sf_3 &:= (Priority, G_1) \\ sf_2 &:= (Duration, G_2) & sf_4 &:= (Duration, G_4) \\ sf_5 &:= (Implementation\ phase\ duration, \perp). \end{aligned}$$

Also, considering s_1 (Example 4), we have:

$$\begin{aligned} \#_{sf_1}(s_1) &= 21 & \#_{sf_3}(s_1) &= 2 \\ \#_{sf_2}(s_1) &= 35 & \#_{sf_4}(s_1) &= 200 \\ \#_{sf_5}(s_1) &= 279 \end{aligned}$$

where s_1 is one of the situations mentioned in 4.

We interpret a nonempty set of situation features, which we call it a *situation feature extraction plan*, as an analog to the schema of tabular data. More formally;

Definition 6 (Situation Feature Extraction Plan). We define a situation feature extraction plan as $SF \subseteq \mathcal{U}_{sf}$ where $SF \neq \emptyset$.

Example 6. A possible situation feature extraction plan for the IT company in Section 2 is as follows:

$$\begin{aligned} SF_{IT} &:= \{(Team\ size, G_3), (Duration, G_2), (Priority, G_1), (Duration, G_4)\} \\ &= \{sf_1, sf_2, sf_3, sf_4\}. \end{aligned}$$

We can map each situation to a data point according to a given situation feature extraction plan. We do that as follows:

Definition 7 (Instance). Let $s \in \mathcal{U}_{situation}$ and $SF \subseteq \mathcal{U}_{sf}$ where $SF \neq \emptyset$. We define the instance $inst_{SF}(s)$ as $inst_{SF}(s) \in SF \rightarrow \mathcal{U}_{val}$ such that $\forall sf \in SF : (inst_{SF}(s))(sf) = \#_{sf}(s)$. We denote the universe of all possible instances as:

$$\mathcal{U}_{instance} = \bigcup_{s \in \mathcal{U}_{situation}} \bigcup_{\substack{SF \subseteq \mathcal{U}_{sf} \\ SF \neq \emptyset}} \{inst_{SF}(s)\}.$$

An instance is a set of pairs where each pair is composed of a situation feature and a value. With a slight abuse of notation, we define $values(sf) = values(at)$ where $sf = (at, G)$ is a situation feature.

Example 7. Considering SF_{IT} from Example 5 and the situations from Example 4. We have:

$$\begin{aligned} inst_{SF_{IT}}(s_1) &= \{((Team\ size, G_3), 21), ((Duration, G_2), 35), ((Priority, G_1), 2), \\ &\quad ((Duration, G_4), 200)\} = \{(sf_1, 21), (sf_2, 35), (sf_3, 2), (sf_4, 200)\} \\ inst_{SF_{IT}}(s_2) &= \{((Team\ size, G_3), 33), ((Duration, G_2), 63), ((Priority, G_1), 1), \\ &\quad ((Duration, G_4), 226)\} = \{(sf_1, 33), (sf_2, 63), (sf_3, 1), (sf_4, 226)\} \\ inst_{SF_{IT}}(s_3) &= \{((Team\ size, G_3), 33), ((Duration, G_2), 63), ((Priority, G_1), 1), \\ &\quad ((Duration, G_4), 62)\} = \{(sf_1, 33), (sf_2, 63), (sf_3, 1), (sf_4, 62)\} \end{aligned}$$

Given a situation feature extraction plan \mathbf{SF} , we consider one of its situation features as the class situation feature, denoted as csf and $\mathbf{SF} \setminus \{csf\}$ as descriptive situation features. Given $\mathbf{SF} \subseteq \mathcal{U}_{sf}$, $csf \in \mathbf{SF}$ where $csf = (at, G)$, and an event log L , we can generate a class situation feature sensitive tabular data-set. We call such a tabular data set a *situation feature table*. To do that, we first generate $S_{L,G}$ and then we generate the situation feature table which is the bag of instances derived from the situations in $S_{L,G}$, regarding \mathbf{SF} . Note that choosing $S_{L,G}$ such that G is the same group as the one in class situation feature (where we have $csf = (at, G)$), ensures the sensitivity of the extracted data to the class situation feature. More formally;

Definition 8 (Situation Feature Table). Let $L \in \mathcal{L}$ be an event log, $\mathbf{SF} \subseteq \mathcal{U}_{sf}$ a situation feature extraction plan, and $csf = (at, G) \in \mathbf{SF}$. We define a situation feature table $T_{L,\mathbf{SF},(at,G)}$ (or equivalently $T_{L,\mathbf{SF},csf}$) as follows:

$$T_{L,\mathbf{SF},(at,G)} = [inst_{\mathbf{SF}}(s) | s \in S_{L,G}].$$

Note that if $csf = (at, G)$ where $G \in \mathcal{G}$, then the situation feature table $T_{L,\mathbf{SF},csf}$ includes the instances derived from the situations in G -based situation subset $S_{L,G}$. However, if $G = \perp$, then it includes the situations derived from the situations in trace-based situation subset $S_{L,\perp}$.

Example 8. Based on Example 7 we have

$$T_{L_{IT},\mathbf{SF}_{IT},(Duration,G_4)} = [inst_{\mathbf{SF}_{IT}}(s_1), inst_{\mathbf{SF}_{IT}}(s_2), inst_{\mathbf{SF}_{IT}}(s_3)].$$

Note that in this example, the class situation feature is $csf = sf_4 = (Duration, G_4)$. Another way to present $T_{L_{IT},\mathbf{SF}_{IT},(Duration,G_4)}$ is as follows:

$sf_1 = (Team\ size, G_3)$	$sf_2 = (Duration, G_2)$	$sf_3 = (Priority, G_1)$	$sf_4 = (Duration, G_4)$
21	35	2	200
33	63	1	226
33	63	1	117

In this table, the first row is corresponding to the $inst_{\mathbf{SF}_{IT}}(s_1)$, the second row is corresponding to the $inst_{\mathbf{SF}_{IT}}(s_2)$, and the third row is corresponding to the $inst_{\mathbf{SF}_{IT}}(s_3)$.

4.2 Structural Equation Model

A structural equation model is a data generating model in the form of a set of equations. Each equation encodes how the value of one of the situation features is determined by the value of other situation features. It is worth noting that these equations are a way to determine how the observational and the interventional distributions are generated and should not be considered as normal equations. More formally³;

Definition 9 (Structural Equation Model (SEM)). Let $T_{L,\mathbf{SF},csf}$ be a situation feature table, in which $L \in \mathcal{L}$, $\mathbf{SF} \subseteq \mathcal{U}_{sf}$, and $csf \in \mathbf{SF}$. The SEM of $T_{L,\mathbf{SF},csf}$ is defined as $\mathcal{EQ} \in \mathbf{SF} \rightarrow Expr(\mathbf{SF})$ where for each $sf \in \mathbf{SF}$, $Expr(sf)$ is an expression of the situation features in $\mathbf{SF} \setminus \{sf\}$ and possibly some noise N_{sf} . It is needed that the noise distributions N_{sf} of $sf \in \mathbf{SF}$ be mutually independent.

³ Definition 9 and 11 are based on [14].

We need \mathbf{SF} to be *causal sufficient*, which means \mathbf{SF} includes all relevant situation features. Based on Definition 9, given a SEM \mathcal{EQ} over the \mathbf{SF} of a situation feature table $T_{L,\mathbf{SF},csf}$, for each $sf \in \mathbf{SF}$, the right side of expression $sf = Expr(\mathbf{SF})$ in \mathcal{EQ} does not include sf .

Given \mathcal{EQ} over the \mathbf{SF} of a situation feature table $T_{L,\mathbf{SF},csf}$, the *parents* of the $sf \in \mathbf{SF}$ is the set of situation features that appear in the right side of expression $\mathcal{EQ}(sf)$. The set of parents of a situation feature includes those situation features with a direct causal effect on it.

Example 9. A possible SEM for the situation feature table mentioned in Example 8 is as follows:

$$\begin{aligned} (Priority, G_1) &= N_{(Priority, G_1)} & N_{(Priority, G_1)} &\sim Uniform(1, 3) \\ (Team\ size, G_3) &= 10(Priority, G_1) + N_{(Team\ size, G_3)} & N_{(Team\ size, G_3)} &\sim Uniform(1, 15) \\ (Duration, G_2) &= 2(team\ size, G_3) + N_{(Duration, G_2)} & N_{(Duration, G_2)} &\sim Uniform(-5, 5) \\ (Duration, G_4) &= 2(Duration, G_2) \times (Priority, G_1) + N_{(Duration, G_4)} & N_{(Duration, G_4)} &\sim Uniform(-100, 100) \\ & & & + (team\ size, G_3) + N_{(Duration, G_4)} \end{aligned}$$

The structure of the causal relationships between the situation features in a SEM can be encoded as a directed acyclic graph, which is called *causal structure*. Given a SEM \mathcal{EQ} on a set of situation features \mathbf{SF} , each vertex in its corresponding causal structure is analogous to one of the situation features in \mathbf{SF} . Let $sf_1, sf_2 \in \mathbf{SF}$, there is a directed edge from sf_1 to sf_2 if sf_1 appears in the right side of expression $\mathcal{EQ}(sf_2)$. More formally,

Definition 10 (Causal Structure). Let \mathcal{EQ} be the SEM of the situation feature table $T_{L,\mathbf{SF},csf}$. We define the corresponding causal structure of \mathcal{EQ} as a directed acyclic graph $(\mathbf{U}, \Rightarrow)$ where $\mathbf{U} = \mathbf{SF}$ and $(sf_1, sf_2) \in \Rightarrow$ if $sf_1, sf_2 \in \mathbf{SF}$ and sf_1 appears in the right side of expression $\mathcal{EQ}(sf_2)$.

In the rest of this paper, we use $sf_1 \rightarrow sf_2$ instead of $(sf_1, sf_2) \in \Rightarrow$.

Having a situation feature table $T_{L,\mathbf{SF},csf}$, the structural equation model of its situation features can be provided by a customer who possesses the process domain knowledge or in a data-driven manner.

Example 10. The causal structure of the SEM mentioned in Example 9 is as depicted in Figure 6.

To predict the effect of manipulating one of the situation features on the other situation features, we need to intervene on the SEM by actively setting the value of one (or more) of its situation features to a specific value (or a distribution). Here, we focus on atomic interventions where the intervention is done on just one of the situation features by actively forcing its value to be a specific value.

Definition 11 (Atomic Intervention). Given an SEM \mathcal{EQ} over \mathbf{SF} where $sf \in \mathbf{SF} \setminus \{csf\}$, and $c \in values(sf)$, the SEM \mathcal{EQ}' after the intervention on sf is obtained by replacing $\mathcal{EQ}(sf)$ by $sf = c$ in \mathcal{EQ} .

Note that the corresponding causal structure of \mathcal{EQ}' (after intervention on sf) is obtained from the causal structure of \mathcal{EQ} by removing all the incoming edges to sf [14].

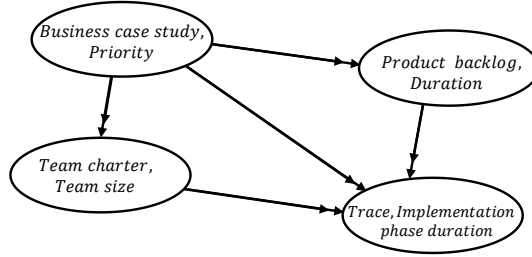


Fig. 6: The causal structure of the SEM mentioned in Example 9.

When we intervene on a situation feature, we just replace the equation of that situation feature in the SEM and the others do not change as causal relationships are autonomous under interventions [14].

Example 11. We can intervene on the SEM introduced in Example 9 by forcing the team size to be 13. For this case, the SEM under the intervention is as follows:

$$\begin{aligned}
 (Priority, G_1) &= N_{(Priority, G_1)} & N_{(Priority, G_1)} &\sim Uniform(1, 3) \\
 (Team\ size, G_3) &= 13 \\
 (Duration, G_2) &= 2(team\ size, G_3) + N_{(Duration, G_2)} & N_{(Duration, G_2)} &\sim Uniform(-5, 5) \\
 (Duration, G_4) &= 2(Duration, G_2) \times (Priority, G_1) + (team\ size, G_3) + N_{(Duration, G_4)} & N_{(Duration, G_4)} &\sim Uniform(-100, 100)
 \end{aligned}$$

5 Approach

Observing a problem in the process, we need to find a set of situation features SF which not only include csf (the situation feature capturing the problem) but also be causal sufficient. The expressiveness of the discovered SEM is highly influenced by SF (even though SEMs, in general, can deal with latent variables). Considering the variety of the possible situation features captured by the event log and the derived ones, finding the proper set SF and also those values of the situation features (or combination of values) that contribute more to the problem is a complicated task and needs plenty of domain knowledge.

We know that correlation does not mean causation. On the other hand, if a situation feature is caused by another situation feature (set of situation features), this implies that there is a correlation between the given situation feature and its parents. We use this simple fact for the automated situation feature recommendation. It is worth noting that the proposed automated situation feature recommendation method is one of the many possible choices. The automated situation feature recommendation method and the SEM discovery process are described in the following:

5.1 Automated situation feature Recommendation

Given an event log $L \in \mathcal{L}$ and the class situation feature $csf = (at, G)$, we declare a nominal situation feature sf as a possible cause of csf if there exists a value $v \in values(sf)$

that appears in big enough portion (at least α where $0 < \alpha \leq 1$) of the situations of $S_{L,G}$ with the undesirable (problematic) result for csf . When sf is a numerical situation feature, we use equal width binning for the values of sf that appear in L where the number of bins, b , is given by the user. We consider sf as a possible cause of csf if there exists a bin of values of sf in L such that the values of that bin appears in more than of α portion of situations in $S_{L,G}$ with the undesirable value for csf . More formally:

Definition 12 (Potential Causal situation feature). Let $L \in \mathcal{L}$ be an event log, $csf = (at, G)$ the class situation feature where $G \in \mathcal{G} \cup \{\perp\}$, α a threshold where $0 < \alpha \leq 1$, and $values(csf)_\downarrow$ denotes the set of undesirable values of csf . We consider a situation feature sf a possible cause of csf if one of the following two conditions holds:

If sf is a nominal situation feature:

$$\exists_{v \in values(sf)} \frac{|\{s \in S_{L,G} | \#_{sf}(s) = v \wedge \#_{csf}(s) \in values(csf)_\downarrow\}|}{|\{s \in S_{L,G} | \#_{csf}(s) \in values(csf)_\downarrow\}|} \geq \alpha.$$

If sf is a numerical situation feature:

$$\exists_{0 \leq i \leq b-1} \frac{|\{s \in S_{L,G} | \#_{sf}(s) \in [\frac{i(v_{max}-v_{min}+1)}{b}, \frac{(i+1)(v_{max}-v_{min}+1)}{b}) \wedge \#_{csf}(s) \in values(csf)_\downarrow\}|}{|\{s \in S_{L,G} | \#_{csf}(s) \in values(csf)_\downarrow\}|} \geq \alpha$$

where $b \in \mathbb{N}$ denotes the number of bins, v_{max} is the maximum value and v_{min} is the minimum value for sf in L .

Moreover, we define the set of all potential causes of csf as the set of all $sf \in \mathcal{U}_{sf}$ for which one of the above inequalities holds.

We present the set of the potential causes to the user as a set of tuples (sf, v) where $sf \in \mathcal{U}_{sf}$ and $v \in values(sf)$ (if sf is a numerical situation feature, then v is the lower bound of the bin) in the descending order regarding the portion of the situations of $S_{L,G}$ with the undesirable result that has value v for sf (has a value in the bin with the lower bound v). This way, the first tuples in the order are those values (lower bound of bins of values) of those situation features that intervention on them may have (potentially) the most effect on the value of the class situation feature.

5.2 SEM Inference

Here we show how to infer the SEM of a given situation feature table in two steps:

- The first step is *causal structure discovery*, which involves discovering its causal structure of the situation feature table. This causal structure encodes the existence and the direction of the causal relationships among the situation features in the situation extraction plan of the given situation feature table.
- The second step is *causal strength estimation*, which involves estimating a set of equations describing how each situation feature is influenced by its immediate causes. Using this information we can generate the SEM of the given situation feature table.

In the sequel, we describe these two steps.

Causal Structure Discovery. The causal structure of the situation features in a given situation feature table can be determined by an expert who possesses the domain knowledge about the underlying process and the causal relationships between its features. But having access to such knowledge is quite rare. Hence, we support discovering the causal structure in a data-driven manner.

Several search algorithms have been proposed in the literature (e.g., [3,20,12]). The input of a search algorithm is observational data in the form of a situation feature table (and possibly knowledge) and its output is a graphical object that represents a set of causal structures that cannot be distinguished by the algorithm. One of these graphical objects is *Partial Ancestral Graph (PAG)* introduced in [23].

A PAG is a graph whose vertex set is $V = SF$ but has different edge types, including $\rightarrow, \leftrightarrow, \bullet\rightarrow, \bullet\bullet$. Similar to \rightarrow , we use infix notation for $\rightarrow, \leftrightarrow, \bullet\rightarrow, \bullet\bullet$. Each edge type has a specific meaning. Let $sf_1, sf_2 \in V$. The semantics of different edge types in a PAG are as follows:

- $sf_1 \rightarrow sf_2$ indicates that sf_1 is a direct cause of sf_2 .
- $sf_1 \leftrightarrow sf_2$ means that neither sf_1 nor sf_2 is an ancestor of the other one, even though they are probabilistically dependent (i.e., sf_1 and sf_2 are both caused by one or more hidden confounders).
- $sf_1 \bullet\rightarrow sf_2$ means sf_2 is not a direct cause of sf_1 .
- $sf_1 \bullet\bullet sf_2$ indicates that there is a relationship between sf_1 and sf_2 , but nothing is known about its direction.

The formal definition of a PAG is as follows [23]:

Definition 13 (Partial Ancestral Graph (PAG)). A PAG is a tuple $(V, \rightarrow, \leftrightarrow, \bullet\rightarrow, \bullet\bullet)$ in which $V = SF$ and $\rightarrow, \leftrightarrow, \bullet\rightarrow, \bullet\bullet \subseteq V \times V$ such that $\rightarrow, \leftrightarrow, \bullet\rightarrow$, and $\bullet\bullet$ are mutually disjoint.

The discovered PAG by the search algorithm represents a class of causal structures that satisfies the conditional independence relationships discovered in the situation feature table and ideally, includes its true causal structure.

Example 12. Two possible PAGs for the SEM mentioned in Example 9 are shown in Figure 7.

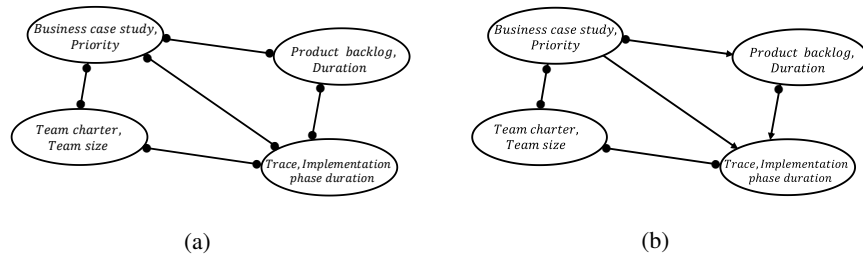


Fig. 7: Two possible PAGs for the SEM mentioned in Example 9.

Now, it is needed to modify the discovered PAG to a compatible causal structure. To transform the output PAG to a compatible causal structure, which represents the causal structure of the situation features in the situation feature table, domain knowledge of the process and common sense can be used. This information can be applied by directly modifying the discovered PAG or by adding them to the search algorithm, as an input, in the form of *required directions* or *forbidden directions* denoted as D_{req} and D_{frb} , respectively. $D_{req}, D_{frb} \subseteq V \times V$ and $D_{req} \cap D_{frb} = \emptyset$. Required directions and forbidden directions influence the discovered PAG as follows:

- If $(sf_1, sf_2) \in D_{req}$, then we have $sf_1 \rightarrow sf_2$ or $sf_1 \bullet \rightarrow sf_2$ in the output PAG.
- If $(sf_1, sf_2) \in D_{frb}$, then in the discovered PAG it should not be the case that $sf_1 \rightarrow sf_2$.

We assume no hidden common confounder exists, so we expect that in the PAG, relation \leftrightarrow be empty⁴. We can define the compatibility of a causal structure with a PAG as follows:

Definition 14 (Compatibility of a Causal Structure With a Given PAG). *Given a PAG $(V, \rightarrow, \leftrightarrow, \bullet \rightarrow, \bullet \bullet \rightarrow)$ in which $\leftrightarrow = \emptyset$, we say a causal structure (U, \Rightarrow) is compatible with the given PAG if $V = U$, $(sf_1 \rightarrow sf_2 \vee sf_1 \bullet \rightarrow sf_2) \implies sf_1 \Rightarrow sf_2$, and $sf_1 \bullet \bullet \rightarrow sf_2 \implies (sf_1 \Rightarrow sf_2 \oplus sf_2 \Rightarrow sf_1)$, where \oplus is the XOR operation and $sf_1, sf_2 \in V$.*

Example 13. The causal structure shown in Figure 6 is compatible with both PAGs demonstrated in Figure 7.

Causal Strength Estimation. The final step of discovering the causal model is estimating the strength of each direct causal effect using the observed data. Suppose D is the causal structure of a situation feature table $T_{L,SF,csf}$. As D is a directed acyclic graph, we can sort its nodes in a topological order γ . Now, we can statistically model each situation feature as a function of the noise terms of those situation features that appear earlier in the topological order γ of D . In other words, $sf = f((N_{sf'})_{sf':\gamma(sf') \leq \gamma(sf)})$ [14]. The set of these functions, for all $sf \in SF$, is the SEM of SF .

Finally, we want to answer questions about the effect of an intervention on any of the situation features on the class situation feature. We can do the intervention as described in Definition 11. The resulting SEM (after intervention) demonstrates the effect of the intervention on the situation features.

Note that, in a given causal structure of a situation feature table $T_{L,SF,csf}$, there is no directed path between $sf \in SF$ and csf , they are independent and consequently, intervening on sf by forcing $sf = c$ has no effect on csf .

6 Experimental Results

We have implemented the proposed approach as a plugin in ProM which is available in the nightly-build of ProM under the name *Root-cause Analysis Using Structural Equation Model*. ProM is an open-source and extensible platform for process mining [21].

⁴ If $\leftrightarrow \neq \emptyset$, the user can restart the procedure after adding some more situation features names to the situation feature table.

The inputs of the implemented plugin are an event log, the Petri-net model of the process, and, the conformance checking results of replaying the given event log on the given model. In the rest of this section, first, we mention some of the implementation details and design choices that we have made and then we present the results of applying the plugin on a synthetic and a real event log.

6.1 Implementation Notes

In the implemented plugin, we first enrich the event log by adding some attributes. Some of the features that can be extracted (besides the ones that have been explicitly mentioned in the event log) from the event log using the implemented plugin are as follows:

- Time perspective: timestamp, activity duration, trace duration, trace delay, sub-model duration.
- Control-flow perspective: next activity, previous activity.
- Conformance perspective: deviation, number of log moves, number of model moves.
- Resource organization perspective: resource, role, group.
- Aggregated features:
 - Process-level: the number of waiting customers, workload.
 - Trace-level: average service time, average waiting time.
 - Event-level: number of active events with a specific activity name, number of waiting events with a specific activity name.
 - Resource-level: average service time, average waiting time

As the second step, the user needs to specify *csf* and *SF*. The user can specify *SF* by manually selecting the proper set of situation features or use the implemented situation feature recommendation method on a predefined set of situation features (for example all the situation features) to identify the relevant set of situation features to *csf*. If *SF* includes an aggregated feature, we also compute the values of the aggregated feature. The third step involves extracting situation feature table from the event log. According to the selected *SF* and *csf* the proper situation subset of the event log is generated and the situation feature table is extracted. Then we infer the causal structure of the situation feature table. For this goal, we use the Greedy Fast Causal Inference (GFCI) algorithm [12] which is a hybrid search algorithm. The inputs of GFCI algorithm are the situation feature table and possibly background knowledge. The output of GFCI algorithm is a PAG with the highest score on the input situation feature table. In [12], it has been shown that under the large-sample limit, each edge in the PAG computed by GFCI is correct if some assumptions hold. Also, the authors of [12] using empirical results on simulated data have shown that GFCI has the highest accuracy among several other search algorithms. In the implemented plugin, we have used the Tetrad [19] implementation of the GFCI algorithm. To use GFCI algorithm, we need to set several parameters. We have used the following settings for the parameters of the GFCI algorithm in the experiments: cutoff for p-values = 0.05, maximum path length = -1, maximum degree = -1, and penalty discount = 2.

In the implemented plugin, we have assumed linear dependencies among the situation features and additive noise when dealing with continuous data. In this case, given

a SEM \mathcal{EQ} over \mathcal{SF} , we can encode \mathcal{EQ} as a weighted graph. This weighted graph is generated from the corresponding causal structure of \mathcal{EQ} by considering the coefficient of sf_2 in $\mathcal{EQ}(sf_1)$ as the weight of the edge from sf_2 to sf_1 . Using this graphical representation of a SEM, to estimate the magnitude of the effect of sf on the csf , we can simply sum the weights of all directed paths from sf to csf , where the weight of a path is equal to the multiplication of the weights of its edges.

6.2 Synthetic Event Log

For the synthetic data, we use the IT company example in Section 2. The situation feature extraction plan is:

$$\{(Team\ size, G_3), (Duration, G_2), (Priority, G_1), (Implementation\ phase\ duration, \perp)\}.$$

where the class situation feature is $(Implementation\ phase\ duration, \perp)$. We assume that the true causal structure of the data is as depicted in Figure 5.

To generate an event log, we first created the Petri-net model of the process as shown in 2 using CPN Tools [17]. Then, using the created model, we generated an event log with 1000 traces. We have enriched the event log by adding *Implementation phase duration* attribute to the traces. This attribute indicates the duration of the sub-model including “development” and “test” transitions in person-day. When generating the log, we have assumed that the true SEM of the process is as follows:

$$\begin{aligned} (Complexity, \perp) &= N_{(Complexity, \perp)} & N_{(Complexity, \perp)} &\sim Uniform(1, 10) \\ (Priority, G_1) &= N_{(Priority, G_1)} & N_{(Priority, G_1)} &\sim Uniform(1, 3) \\ (Duration, G_2) &= 10(Complexity, \perp) + N_{(Team\ size, G_3)} & N_{(Duration, G_2)} &\sim Uniform(-2, 4) \\ (Team\ size, G_3) &= 5(Complexity, \perp) + 3(Priority, G_1) + N_{(Team\ size, G_3)} & N_{(Team\ size, G_3)} &\sim Uniform(-1, 2) \\ (Implementation\ phase\ duration, \perp) &= 50(Complexity, \perp) + & N_{(Implementation\ phase\ duration, \perp)} &\sim Uniform(10, 20) \\ &5(Team\ size, G_3) + N_{(Implementation\ phase\ duration, \perp)} \end{aligned}$$

The summary of the generated event log and its trace variants (generated by ProM) are shown in Figure 8.

Generating situation feature table. We generate a situation feature table using the mentioned situation feature extraction plan. A snapshot of the generated situation feature using the implemented plugin is shown in Figure 9. In this figure, the class situation feature is colored in pink and the descriptive situation features are colored gray.

SEM inference. Applying the implemented plugin on the situation feature extracted from the event log, the PAG depicted in Figure 10a was discovered. Even though the discovered PAG does a good job regarding discovering the potential causal relationship, it does not say much about the direction of them. Here the customer may guess that another influential attribute might exist that acts as a confounder. Considering $(Complexity, \perp)$ as another descriptive situation feature, then the discovered PAG by the implemented plugin would be as the one in Figure 10b. This PAG is more accurate and includes the true causal structure of the situation feature table. We have assumed that the complexity of a project is a feature that is not recorded in the event log. The customer, may assume (based on domain knowledge) that the duration of “product backlog” is longer in more complex projects and assign to the complexity of a project the floor

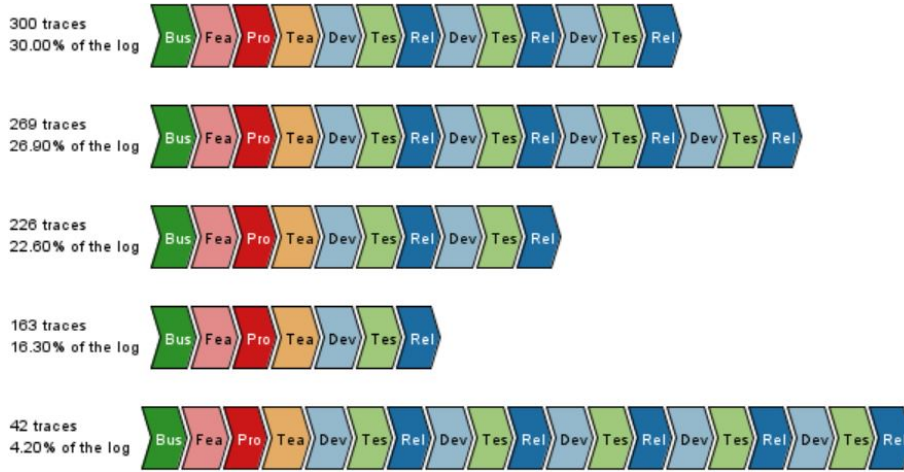


Fig. 8: The trace variants in the synthetic event log.

Team size : Team charter	Duration person day : Pro...	Priority : Business case d...	Implementation phase d...
20	32	1	266
34	49	3	432
15	12	3	138
42	71	2	577
20	30	2	265
8	12	1	100
46	72	3	590
45	82	2	641
30	52	2	413
54	91	3	740
28	41	3	357
18	20	3	200

Fig. 9: A snapshot of the situation feature table generated for the synthetic event log.

of the value of $(Duration, G_2)$ divided by 10. Now, using domain knowledge and the chronological order of transitions, we can turn the discovered PAG into the causal structure depicted in Figure 10c. After estimating the strength of the causal relationships, we obtain the SEM shown in Figure 10d.

Moreover, we can have the inferred SEM in text format. In this case, the output would be as Shown in Figure 11.

By comparing the estimated coefficients of situation features names in the output of the plugin (and equivalently the weights of the edges in Figure 10d), and those in the equations of the true SEM of the data, we can see that the estimated and real strengths of causal relationships are quite close.

To investigate the effect of an intervention on any of the situation features on the class situation feature, we can find the equation capturing the effect of intervention by simply clicking on its corresponding node in the causal structure. For example, if we

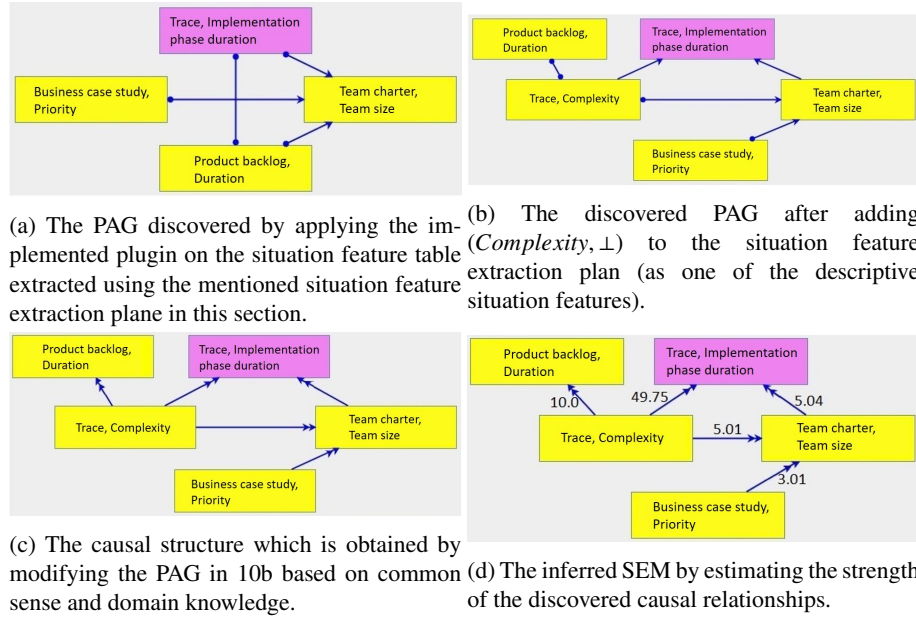


Fig. 10: The PAC, causal structure and the SEM discovered using implemented plugin for the synthetic event log.

$$\begin{aligned}
 \text{Team size} : \text{Team charter} &= \text{Complexity} * 5.01 + \text{Priority} : \text{Business case development} * 3.01 + \text{noise} \\
 \text{Duration person day} : \text{Product backlog} &= \text{Complexity} * 10.0 + \text{noise} \\
 \text{Priority} : \text{Business case development} &= \text{noise} \\
 \text{Implementation phase duration person day} &= \text{Team size} : \text{Team charter} * 5.04 + \text{Complexity} * 49.75 + \text{noise} \\
 \text{Complexity} &= \text{noise}
 \end{aligned}$$

Fig. 11: The discovered SEM from the situation feature table extracted from the synthetic event log after adding (*Complexity*, \perp) to the situation feature extraction plan in the text format.

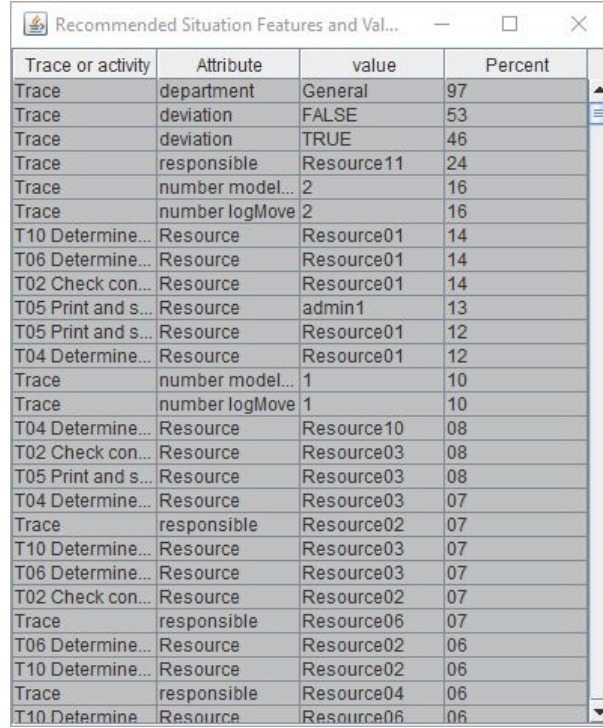
click on the corresponding node of (*Team size*, G_3), we have

$$(\text{Implementation phase duration}, \perp) = 75.0004 \times (\text{Complexity}, \perp) + \text{noise}.$$

This equation means that by enforcing the complexity of a project to be one unit more complex, then we expect that its implementation phase takes approximately 75 more person-days (assuming that the complexity of a project is actionable). As another example, equation $(\text{Implementation phase duration}, \perp) = 0.0 \times (\text{Duration}, G_2)$ shows the estimated effect of intervention on (*Duration*, G_2). We can interpret this equation as “intervention on (*Duration*, G_2) has no effect on (*Implementation phase duration*, \perp)”.

6.3 Real Event Log

We have used the implemented plugin also on several real-life event logs. In this subsection we analyse *receipt phase of an environmental permit application process (WABO)*



Trace or activity	Attribute	value	Percent
Trace	department	General	97
Trace	deviation	FALSE	53
Trace	deviation	TRUE	46
Trace	responsible	Resource11	24
Trace	number model...	2	16
Trace	number logMove	2	16
T10 Determine...	Resource	Resource01	14
T06 Determine...	Resource	Resource01	14
T02 Check con...	Resource	Resource01	14
T05 Print and s...	Resource	admin1	13
T05 Print and s...	Resource	Resource01	12
T04 Determine...	Resource	Resource01	12
Trace	number model...	1	10
Trace	number logMove	1	10
T04 Determine...	Resource	Resource10	08
T02 Check con...	Resource	Resource03	08
T05 Print and s...	Resource	Resource03	08
T04 Determine...	Resource	Resource03	07
Trace	responsible	Resource02	07
T10 Determine...	Resource	Resource03	07
T06 Determine...	Resource	Resource03	07
T02 Check con...	Resource	Resource02	07
Trace	responsible	Resource06	07
T06 Determine...	Resource	Resource02	06
T10 Determine...	Resource	Resource02	06
Trace	responsible	Resource04	06
T10 Determine...	Resource	Resource06	06

Fig. 12: The first 27 recommended situation features and values by applying our new situation feature recommendation method.

CoSeLoG project [2] event log (receipt log for short). The receipt event log includes 1434 traces and 8577 events. The maximum duration of traces in the receipt event log is almost 276 days while and the average duration of traces is almost 3 minutes. We consider those traces that took longer than 1 percent of the maximum trace duration as delayed. Thus, the class situation feature is (*trace delay*, \perp) and “trace delay” is one of the trace-level attributes that has been used to enrich the receipt event log. The length of the time window in this experiment has been set to one day.

Situation feature extraction plan selection. In this case study, we first use the situation feature recommendation in which we set the number of bins to 100. We use the set including features related to the resources of events, the duration of events, process workload, and some of the trace features (such as deviation, number of model moves, number of log moves, responsible, and department) as the initial set of features. In Figure 12, one can see the first 27 recommended situation features and values. The last column shows in which percent of the situations with the undesirable result, the situation feature has the mentioned value.

We set the α threshold to 0.05%. As a result, the situation feature extraction plan includes some of the resources for some of the events, responsible, deviation, num-

Process workL.	responsible	Trace Delay	Resource : T02	Resource : T05	Resource : T06	deviation	Resource : T04	number logMove	Resource : T10
16	Resource21	delayed	Resource21	NOT SET	NOT SET	true	NOT SET	1	NOT SET
11	Resource04	onTime	Resource30	NOT SET	Resource30	true	NOT SET	5	Resource30
11	Resource11	onTime	Resource03	Resource03	Resource03	false	Resource03	0	Resource03
13	Resource11	onTime	Resource03	Resource03	Resource03	false	Resource03	0	Resource03
13	Resource11	onTime	Resource03	Resource03	Resource03	true	Resource03	4	Resource03
8	Resource06	onTime	Resource06	Resource06	Resource06	false	Resource06	0	Resource06
15	Resource11	onTime	Resource03	Resource03	Resource03	true	Resource03	2	Resource03
11	admin3	onTime	NOT SET	NOT SET	NOT SET	true	NOT SET	0	NOT SET
8	Resource06	onTime	Resource06	Resource06	Resource06	false	Resource06	0	Resource06
13	Resource22	onTime	Resource22	Resource22	Resource22	false	Resource22	0	Resource22
18	Resource13	onTime	Resource13	Resource13	Resource13	false	Resource13	0	Resource13

Fig. 13: A snapshot of the situation feature table extracted from the receipt event log based on the selected situation feature extraction plan.

ber of log moves, number of model moves, and process workload. We have ignored some of the situation features recommended by the algorithm. For example, we have ignored (*department*, \perp) as the value assigned to this situation feature in almost all of the traces is “General”, so it does not contain any information about the class situation feature. Also, we have ignored (*deviation*, \perp) as each of its assigned values appears in almost the same portion of the situations with the undesirable result. We have ignored (*Number logMove*, \perp) as it has the same value as (*Number modelMove*, \perp) in all of the generated instances. For the sake of simplicity, we use the following abbreviation for activity names in the rest of this section.

- “T02” instead of “T02 Check confirmation of receipt”.
- “T04” instead of “T04 Determine confirmation of receipt”.
- “T05” instead of “T05 Print and send confirmation of receipt”.
- “T06” instead of “T06 Determine necessity of stop advice”.
- “T10” instead of “T10 Determine necessity to stop indication”.

The situation feature extraction plan in this example is as follows:

$$\{(G_{T02}, Resource), (G_{T04}, Resource), (G_{T05}, Resource), (G_{T06}, Resource), (G_{T10}, Resource), (\perp, Trace\ delay), (\perp, Process\ workload), (\perp, Responsible), (\perp, Deviation), (\perp, Number\ logMove)\}$$

where the class situation feature is (\perp , *Trace delay*) and

$$\begin{aligned} (G_{T02} &:= group(actName, \{“T02”\})) & (G_{T05} &:= group(actName, \{“T05”\})) \\ (G_{T04} &:= group(actName, \{“T04”\})) & (G_{T06} &:= group(actName, \{“T06”\})) \\ (G_{T10} &:= group(actName, \{“T10”\})). \end{aligned}$$

Generating situation feature table. Using the above situation feature extraction plan, we generate a situation feature table. A snapshot of the generated situation feature table is shown in Figure 13.

SEM inference. The discovered PAG from the extracted situation feature is as shown in Figure 14. Using the temporal ordering of the activities (in this process, *T02* happens before *T04*, *T04* happens before *T05*, *T05* happens before *T06*, *T06* happens before *T10* in all the traces) and common sense (the choice of the resource of an activity does

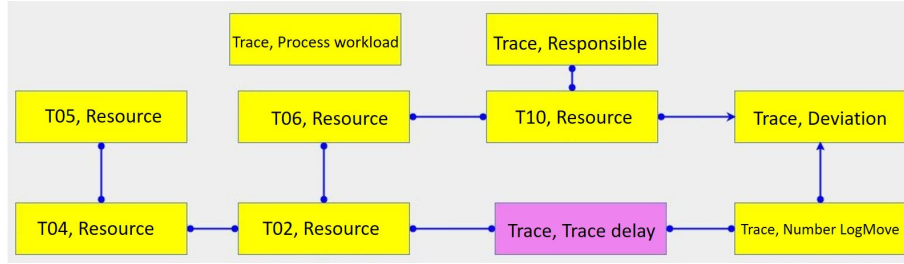


Fig. 14: The PAG of the situation feature table extracted from the receipt event log, discovered by the implemented plugin.

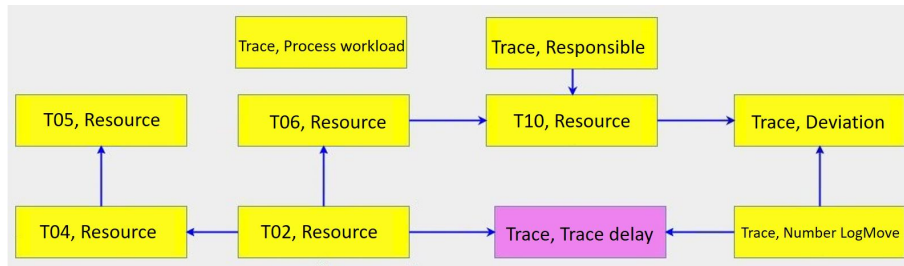


Fig. 15: The causal structure obtained by modifying the PAG in Figure 14 using common sense and domain knowledge.

not affect the choice of the resource of another activity that happened before) we are able to convert the PAG in Figure 14 into the causal structure shown in Figure 15.

According to this causal structure only $(G_{T02}, Resource)$ and $(\perp, Number\ logMove)$ have a causal effect on $(\perp, Trace\ delay)$ and there is no causal relationship between other situation features and $(\perp, Trace\ delay)$. Some of the other causal relationships encoded in Figure 15 are as follows:

- The choice of $(G_{T02}, Resource)$ directly influences the choice of $(G_{T04}, Resource)$ and $(G_{T06}, Resource)$.
- The value of $(G_{T02}, Resource)$ indirectly influences the value of $(G_{T05}, Resource)$, $(G_{T10}, Resource)$.
- $(\perp, responsible)$ and $(G_{T06}, Resource)$ causally influence $(G_{T10}, Resource)$.
- $(\perp, Number\ logMove)$ and $(G_{T10}, Resource)$ directly influence $(\perp, Deviation)$.
- There is no causal relationships between any of the situation features and $(\perp, Process\ workload)$.

After applying the estimation step, we can see the interventional distributions of $(\perp, Trace\ delay)$. We can see the effect of an intervention on any of the situation features by clicking on its corresponding node. For example, we can see that the probability of $(\perp, Trace\ delay) = delayed$ by forcing $(G_{T02}, Resource) = Resource14$ is almost 0.256. The predicted interventional distributions resulting from enforcing $(G_{T02}, Resource)$ to be assigned each of its possible values is shown in Figure 16.

Resource : T02 Chec...	delayed	onTime
Resource22	0.153846153846153...	0.8461538461538461
Resource23	0.25	0.75
Resource24	0.9354838709677419	0.064516129032258...
Resource25	0.3333333333333333	0.6666666666666666
Resource26	0.0	1.0
Resource28	1.0	0.0
Resource29	0.5	0.5
Resource30	0.0	1.0
Resource31	0.0	1.0
Resource32	1.0	0.0
Resource10	0.7049180327868853	0.295081967213114...
Resource33	1.0	0.0
Resource11	0.181818181818181...	0.8181818181818182
Resource34	0.5	0.5
Resource12	0.2727272727272727	0.7272727272727273
Resource35	0.0	1.0
Resource13	0.2857142857142857	0.7142857142857143
Resource36	0.0	1.0
Resource14	0.2558139534883721	0.7441860465116279
Resource38	0.0	1.0
Resource15	0.59375	0.40625
Resource39	1.0	0.0
Resource16	0.027777777777777...	0.9722222222222222
admin1	1.0	0.0
Resource17	0.25	0.75
admin2	0.4444444444444444	0.5555555555555556
Resource18	0.4222222222222222	0.5777777777777778

Fig. 16: The predicted interventional distributions resulting from intervention on ($G_{T02}, Resource$) by enforcing different values.

It is worth noting that as the data in this case study is categorical, the inferred SEM includes several equations where the right hand side includes many terms. Each term shows the distribution of each one of the situation feature values condition on one of the values of one of its direct causes. So, it is not possible to completely present the inferred SEM like in the continuous case.

7 Conclusion

Distinguishing causal from mere correlational relationships among the process features is a vital task when investigating the root causes of performance and/or conformance problems in a company. The best way to identify the causal relationships is using randomized experiments. However, this requires implementing process changes to see their effect. As applying randomized experiments is usually quite expensive (if not impossible) in the processes, we propose a method for root cause analysis based on the theory of causality which uses a mixture of data analysis and domain knowledge. The stakeholders can use this framework to incorporate both domain knowledge and potential statistically supported causal effects to find the SEM of the features and indicators of the process. Moreover, this method helps stakeholders investigating the effect of an

intervention on the process. This information can be used to design and order the re-engineering steps.

The validity of a discovered structural equation model (and any other machine learning technique) is highly influenced by the set of features that have been used for data extraction and structural equation model discovery. However, the complex and dynamic inter-dependencies in processes makes the task of selecting the set of features with a potential causal effect on the observed problem in the process a challenging task. So, we have proposed a simple yet intuitive and effective feature recommendation method in this paper. The proposed method provides the user not just the set of features with the possible causal effect on the class situation feature but also those values of the features that potentially have the largest contribution to the observed problem in the process. It is worth noting that this was missing in the previous work on the causal structure model discovery of a process.

As future work, we would like to learn more process-related features that go beyond individual cases. For example, bottlenecks are caused by competing cases or a shortage of resources. Also, notions such as blocking, batching, and overtaking are not captured well. We would also like to make the diagnostics more understandable. This requires mapping diagnoses related to features back onto the process model and even log. Finally, we would like to enhance simulation models with SEM-based rules.

Acknowledgement

We thank the Alexander von Humboldt (AvH) Stiftung for supporting our research.

References

1. Bozorgi, Z.D., Teinmaa, I., Dumas, M., La Rosa, M., Polyvyanyy, A.: Process mining meets causal machine learning: Discovering causal rules from event logs. In: 2020 2nd International Conference on Process Mining (ICPM). pp. 129–136. IEEE (2020)
2. Buijs, J.: Receipt phase of an environmental permit application process ('wabo'), coselog project. Eindhoven University of Technology (2014), <http://dx.doi.org/10.4121/uuid:a07386a5-7be3-4367-9535-70bc9e77dbe6>
3. Chickering, D.M.: Optimal structure identification with greedy search. *Journal of machine learning research* **3**(Nov), 507–554 (2002)
4. Gupta, N., Anand, K., Sureka, A.: Pariket: Mining business process logs for root cause analysis of anomalous incidents. In: *Proceedings of Databases in Networked Information Systems - 10th International Workshop*. vol. 8999, pp. 244–263. Springer (2015). https://doi.org/10.1007/978-3-319-16313-0_19
5. Hompes, B.F.A., Maaradji, A., Rosa, M.L., Dumas, M., Buijs, J.C.A.M., van der Aalst, W.M.P.: Discovering causal factors explaining business process performance variation. In: *Proceedings of Advanced Information Systems Engineering*. vol. 10253, pp. 177–192. Springer (2017). https://doi.org/10.1007/978-3-319-59536-8_12
6. Lehto, T., Hinkka, M.: Discovering business area effects to process mining analysis using clustering and influence analysis. In: *International Conference on Business Information Systems*. pp. 236–248. Springer (2020)

7. Lehto, T., Hinkka, M., Hollmén, J.: Focusing business improvements using process mining based influence analysis. In: International Conference on Business Process Management. pp. 177–192. Springer (2016)
8. Lehto, T., Hinkka, M., Hollmén, J., et al.: Focusing business process lead time improvements using influence analysis. In: SIMPDA. pp. 54–67 (2017)
9. de Leoni, M., van der Aalst, W.M., Dees, M.: A general process mining framework for correlating, predicting and clustering dynamic behavior based on event logs. *Inf. Syst.* **56**(C), 235–257 (2016). <https://doi.org/10.1016/j.is.2015.07.003>
10. Mothilal, R.K., Sharma, A., Tan, C.: Explaining machine learning classifiers through diverse counterfactual explanations. In: Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency. pp. 607–617 (2020), <http://arxiv.org/abs/1905.07697>
11. Narendra, T., Agarwal, P., Gupta, M., Dechu, S.: Counterfactual reasoning for process optimization using structural causal models. In: Proceedings of Business Process Management Forum. vol. 360, pp. 91–106. Springer (2019). https://doi.org/10.1007/978-3-030-26643-1_6
12. Ogarrio, J.M., Spirtes, P., Ramsey, J.: A hybrid causal search algorithm for latent variable models. In: Proceedings of Probabilistic Graphical Models - Eighth International Conference. pp. 368–379 (2016), <http://proceedings.mlr.press/v52/ogarrio16.html>
13. Pearl, J.: Causality. Cambridge university press (2009)
14. Peters, J., Janzing, D., Schölkopf, B.: Elements of causal inference: foundations and learning algorithms. MIT press (2017)
15. Qafari, M.S., van der Aalst, W.M.: Root cause analysis in process mining using structural equation models. In: International Conference on Business Process Management. pp. 155–167. Springer (2020)
16. Qafari, M.S., van der Aalst, W.M.: Case level counterfactual reasoning in process mining. arXiv preprint arXiv:2102.13490 (2021)
17. Ratzer, A.V., Wells, L., Lassen, H.M., Laursen, M., Qvortrup, J.F., Stissing, M.S., Westergaard, M., Christensen, S., Jensen, K.: CPN tools for editing, simulating, and analysing coloured petri nets. In: Proceedings of Applications and Theory of Petri Nets. vol. 2679, pp. 450–462. Springer (2003). https://doi.org/10.1007/3-540-44919-1_28
18. Sani, M.F., van der Aalst, W.M.P., Bolt, A., García-Algarra, J.: Subgroup discovery in process mining. In: Proceedings of Business Information Systems. pp. 237–252. Springer (2017). https://doi.org/10.1007/978-3-319-59336-4_17
19. Scheines, R., Spirtes, P., Glymour, C., Meek, C., Richardson, T.: The tetrad project: Constraint based aids to causal model specification. *Multivariate Behavioral Research* **33**(1), 65–117 (1998)
20. Spirtes, P., Glymour, C.N., Scheines, R., Heckerman, D.: Causation, prediction, and search. MIT press (2000)
21. Verbeek, H., Buijs, J., Van Dongen, B., van der Aalst, W.M.: Prom 6: The process mining toolkit. *Proc. of BPM Demonstration Track* **615**, 34–39 (2010)
22. Wang, Y., Liang, D., Charlin, L., Blei, D.M.: The deconfounded recommender: A causal inference approach to recommendation. *CoRR* [abs/1808.06581](https://arxiv.org/abs/1808.06581) (2018), <http://arxiv.org/abs/1808.06581>
23. Zhang, J.: On the completeness of orientation rules for causal discovery in the presence of latent confounders and selection bias. *Artif. Intell.* **172**(16-17), 1873–1896 (2008). <https://doi.org/10.1016/j.artint.2008.08.001>