

Mining and Reasoning on Workflows

Gianluigi Greco, Antonella Guzzo, Giuseppe Manco, and
Domenico Saccà, *Senior Member, IEEE Computer Society*

Abstract—Today's workflow management systems represent a key technological infrastructure for advanced applications that is attracting a growing body of research, mainly focused in developing tools for workflow management, that allow users both to specify the "static" aspects, like preconditions, precedences among activities, and rules for exception handling, and to control its execution by scheduling the activities on the available resources. This paper deals with an aspect of workflows which has so far not received much attention even though it is crucial for the forthcoming scenarios of large scale applications on the Web: Providing facilities for the human system administrator for identifying the choices performed more frequently in the past that had lead to a desired final configuration. In this context, we formalize the problem of discovering the most frequent patterns of executions, i.e., the workflow substructures that have been scheduled more frequently by the system. We attacked the problem by developing two data mining algorithms on the basis of an intuitive and original graph formalization of a workflow schema and its occurrences. The model is used both to prove some intractability results that strongly motivate the use of data mining techniques and to derive interesting structural properties for reducing the search space for frequent patterns. Indeed, the experiments we have carried out show that our algorithms outperform standard data mining algorithms adapted to discover frequent patterns of workflow executions.

Index Terms—Data mining, workflow management.

1 INTRODUCTION

A workflow is a partial or total automation of a business process in which a collection of *activities* must be executed by humans or machines according to certain procedural rules. Modern enterprises increasingly use workflow technology for designing business processes by means of management systems that provide mechanisms for formally specifying the schema of execution, for simulating its evolution under different conditions, for validating and testing whether it behaves as expected, and for evaluating the ability of a service to meet requirements with respect to throughput times, service levels, and resource utilization.

This paper deals with an aspect of workflows which has so far not received much attention even though it is crucial for the forthcoming scenarios of large scale applications on the Web: Providing facilities for the human system administrator to monitor the actual behavior of the workflow system in order to predict the "most probable" workflow executions. Indeed, in real-world cases, the enterprise must perform many choices during workflow execution; some choices may lead to a benefit, others should instead be avoided in the future. Data mining techniques may, obviously, help the administrator by looking at all of the previous instantiations (suitably collected into log files in any commercial system) in

order to extract unexpected and useful knowledge about the process and in order to take the appropriate decisions in the executions of future instances. The discovered knowledge can be profitably used for solving problems such as:

- **Successful Termination Prediction.** Assume that an execution is at a given point in which the administrator has to choose an activity to start from a given set of potential activities. Then, she/he typically wants to know which is the choice performed in the past that had more frequently led to a desired final configuration.
- **Identification of Critical Activities.** In every workflow schema, there are some activities that can be considered *critical* in the sense that they are scheduled by the system in every successful execution. Sometimes, the system administrator may know in advance that a given activity is critical, but it often happens that this knowledge must be inferred by looking at the actual behavior of the system.
- **Failure/Success Characterization.** By analyzing the past experience, a workflow administrator may be interested in knowing which discriminant factors characterize the failure or the success in the executions.
- **Workflow Optimization.** The information collected into the logs of the system can be profitably used to reason on the "optimality" of workflow executions. For instance, the optimality criterion can be fixed with regard to some real-case interesting parameter, such as the quality of the service or the average completion time.

In this paper, we concentrate on the first of the above problems: *Successful Termination Prediction*. We show that a crucial step towards an automatic solution to this problem consists of identifying the blocks of activities, called *patterns*, that have been more frequently scheduled together during the execution by the workflow system. To this aim, we propose two distinct algorithms for *frequent pattern mining*

• G. Greco is with the Department of Mathematics, University of Calabria, Via Bucci 30B, I87036 Rende (CS), Italy. E-mail: ggrec@mat.unical.it.

• A. Guzzo is with the DEIS Department, University of Calabria, Via Bucci 41c, I87036 Rende (CS), Italy. E-mail: guzzo@deis.unical.it.

• G. Manco is with the Institute of High Performance Computing and Networks (ICAR-CNR), Via Bucci 41c, I87036 Rende (CS), Italy. E-mail: manco@icar.cnr.it.

• D. Saccà is with the Institute of High Performance Computing and Networks (ICAR-CNR) and the DEIS Department, University of Calabria, Via Bucci 41c, I87036 Rende (CS), Italy. E-mail: sacca@unical.it.

Manuscript received 29 April 2004; revised 8 August 2004; accepted 20 October 2004; published online 17 Feb. 2005.

For information on obtaining reprints of this article, please send e-mail to: tkde@computer.org, and reference IEEECS Log Number TKDE-0147-0803.

implementing sophisticated techniques which benefit from the peculiarities of the applicative context, thereby extending previous proposals for mining frequent structures in complex domains (such as frequent sequences, trees, and graphs—see, e.g., [1], [2], [3], [4], [5], [6], [7], [8]) to the mining of workflow executions.

1.1 Related Work

This paper is about applying data mining techniques to the area of workflows and, as such, it presents quite an intuitive graph formalization of the main workflow concepts as the basic data structure on which data mining algorithms work. Therefore, the first area of related work is workflow modeling and analysis. Let us preliminarily point out that this paper is not aimed at developing a comprehensive workflow specification; so, even though our workflow model covers basic features required in workflow specification, it contains some simplifying assumptions. For instance, our model does not incorporate compensation or reset activities and assumes acyclicity, i.e., nonrecursive workflows and noniterative executions. Furthermore, our model does not directly support scheduling or verification tasks (see, e.g., [9]) and does not handle transactional properties of processes. The reader interested in the description of advanced features in workflow modeling is referred to [10], which proposes a unifying model for concurrency control and recovery for processes. Other elaborated models are: the *Concurrent Transaction Logic-based model (CTR)* [11], which enables to both describe and reason about workflow by introducing a rich set of constraints [12]; the *state chart model* [13], [14], in which *triggers* are introduced in order to define *ECA* (Event Condition Action) rules for describing transitions among states; the *active object oriented model* [15], in which a workflow is modeled by integrating *ECA* rules with object-oriented concepts; the *Process algebra-based model* [16]; and, the *Petri Nets-based model* [17], which is a formalism having deep formal foundations and is profitably used for investigating different interesting properties for the process, such as liveness and boundness (see, e.g., [18]).

Let us now review some work related to the main topic of this paper, that is, graph mining techniques specialized to handle constraints derived by the structures of workflow schemes and instances. The idea of *mining* execution traces has been already addressed in the context of process discovery (see, e.g., [19]), but the goal there is to use the information collected at runtime for deriving an “a posteriori” schema that can model all logged executions and is well-suited for adapting the system to changing circumstances and removing imperfections in the initial design (see [20] for a survey on this topic).

Instead, in our approach, the workflow schema is the starting point, not the result: A number of executions are analyzed contextually and comparatively on the basis of the schema and with the goal of finding frequent patterns of activities, thus discovering useful knowledge for supporting the decision process in an enterprise. To the best of our knowledge, our work is the first in handling such a problem, that is, a problem of mining graphs with constraints imposed by the structures of workflow schemes and instances. Under this perspective, the techniques we propose must be compared with other efforts paid by the database community for developing algorithms for mining frequent patterns, both in relational databases and in complex domains. Most of these approaches are based on the *antimonotone* property, first exploited in the seminal paper of Agrawal and Srikant [21] that introduces the *Apriori*

algorithm: The idea is to generate the set of candidates of length $k + 1$ by combining in a suitable way the set of frequent patterns of size k and then checking their frequencies. Quite a simple generalization of this method is presented in [1] in order to mine sequential patterns.

A completely different approach has been proposed in [3] and goes by the name of the *FP-growth* method. Essentially, the idea is to mine frequent instances with a top-down approach, i.e., by recursively projecting the database according to the frequent patterns already found and then combining the results of mining the projected databases. The extension to a sequential pattern is the *PrefixSpan* algorithm [4]. A recent attempt for combining such a method with the *Apriori* approach has been done in [22].

As for the problem of mining patterns in complex domains, the discovery of frequent trees in a forest has been tackled in [2], while a first *Apriori*-based algorithm, called *AGM*, for identifying frequent substructures in a given graph data set has been presented in [5]. We stress that this latter task nowadays constitutes a very active and still promising area of research for its interesting applications in web analysis and bioinformatics.

For instance, the level-wise search performed by *AGM* has been adopted and further improved in the *FSG* algorithm [7], in which a smart strategy for labeling the generated subgraphs avoids many computational expensive subgraph isomorphism computations. Moreover, some algorithms based on the projection method have quite recently been proposed as well: *gSpan* [6] discovers all the frequent subgraphs without candidate generation and false positive pruning, whereas *CloseGraph* [8] dramatically reduces the number of unnecessary frequent subgraphs generated by exploiting the notion of *closed* patterns, i.e., patterns which are no proper subgraphs of any other pattern with the same support.

It is clear that such approaches could, in principle, be used to deal with the problem of mining frequent workflow instances after a suitable adaptation for fitting the peculiarities of the specific applicative domain of workflow systems. In fact, one can think of modeling the workflow schema as a graph and the executions of the workflow as a set of subgraphs complying with the graph representing the workflow schema.

However, the adaptation of the above mentioned methods to workflow mining is a challenging task and its results are impractical from both the expressiveness and efficiency viewpoints. Indeed, generation of patterns with such traditional approaches does not benefit from the exploitation of the executions’ constraints imposed by the workflow schema, such as precedences among activities, synchronization, and parallel executions of activities (see, e.g., [23], [24], [25]). In contrast, the algorithms proposed in this paper are novel mining techniques *specialized* to handle constraints derived by the structures of workflow schemes. And, in fact, several experiments, reported in Section 5, confirmed that they outperform traditional data mining algorithms, even though they are suitably reengineered (in our implementations) to work with workflow instances.

We conclude the overview on related work by observing that, in order to model all the details of a workflow system, one viable way is to consider more expressive approaches, such as the multirelational data mining approaches [26]. Nonetheless, in Section 5, we also show that as a consequence of their generality in modeling different domains, they poorly perform if compared with our algorithms specifically designed for the workflow domain.

1.2 Contribution of the Paper

In this paper, we investigate the possibility of exploiting data mining techniques within workflow management contexts by proposing two algorithms for mining frequent workflow patterns of execution. Specifically, our contribution is as follows:

- We model the *Successful Termination Problem* and we provide some intractability results that shed light into the intrinsic difficulty on reasoning over workflows. The in-depth theoretical analysis we provide strongly motivates the use of Data Mining techniques, thus confirming the validity of the approach.
- We define the notion of workflow *pattern* and *weak pattern* of a workflow graph, where the latter is a syntactic restriction of the former. In particular, we prove that weak patterns are well-suited for mining tasks as they can be recognized in a highly parallelizable way and can be easily composed to discover frequent patterns because of their interesting structural properties. Indeed, we show that the space of all connected weak patterns constitutes a lower semilattice with regard to a particular relation precedence (\prec).
- By exploiting properties of weak patterns, we design two algorithms for mining frequent patterns that conform to the workflow specifications:
 - *w*-find, which performs a smart (level-wise) exploration of the lower semilattice, and
 - *c*-find, which mines frequent instances by composing connected components.
- We test *w*-find and *c*-find by evaluating their performance and their scalability. We show that none of the algorithms is the best in absolute terms by also evidencing the discriminant factors. Moreover, we compare these algorithms with existing techniques adapted to our particular domain. Several experiments confirmed the validity and the usefulness of these approaches.

We stress that our approach does not consider cyclic graphs (i.e., recursive workflow schemas and iterated executions) and other aspects of workflows, such as compensation or reset activities. These assumptions have been required by the necessity of starting from a simplified model, yet covering important and typical features required in workflow specification, to take up an interesting and relevant topic that has not been given much attention in the literature so far. In fact, a significant number of technical challenges had to be faced for dealing with even basic features only. However, since there is no conceptual limitation in extending our algorithms for mining frequent instances with regard to more involved workflow models, we believe that this work might stimulate the data mining community in continuing our investigation and in facing some of the challenges we posed here.

1.3 Organization

The paper is organized as follows: Section 2 provides a formal model of workflows and many complexity considerations on such a proposed model. A formalization of the problems of *Successful Termination Prediction* and of mining frequent patterns from workflow schemas is devised in Section 3. The levelwise theory of workflow patterns is presented in Section 4, together with the algorithms *w*-find

and *c*-find. Finally, Section 5 provides an experimental validation of the approach.

2 THE WORKFLOW ABSTRACT MODEL

A significant amount of research has been already done in the specification of mechanisms for process modeling (see, e.g., [27] for an overview of different proposals). The most widely adopted formalism is the *control flow graph*, in which a workflow is represented by a labeled directed graph whose nodes correspond to the tasks to be performed, and whose arcs describe the precedences among them. Moreover, the Workflow Management Coalition¹ has also identified additional controls, such as loops and subworkflows.

In this paper, we do not refer to any particular model proposed in the literature. Rather, we next provide a simple (state based) model that covers most of the important and typical features required in workflow specification. The model will be used for providing, in a rigorous way, both the syntax and the execution semantics. Hence, it will trace the formal framework (whose limitations have been already described in the Introduction) for developing our mining algorithms.

Definition 1. A workflow schema WS is a tuple $\langle A, E, a_0, F, IN, OUT_{min}, OUT_{max} \rangle$, where A is a finite set of activities, $E \subseteq (A - F) \times (A - \{a_0\})$ is an acyclic relation of precedences among activities, $a_0 \in A$ is the starting activity, $F \subseteq A$ is the set of final activities, and IN, OUT_{min} , and OUT_{max} are three functions assigning to each node a natural number ($A \mapsto \mathbb{N}$) as follows:

- $\forall a \in A - \{a_0\}, 0 < IN(a) \leq InDegree(a);$
 - $\forall a \in A - F,$
 $0 < OUT_{min}(a) \leq OUT_{max}(a) \leq OutDegree(a);$
 - $IN(a_0) = 0$ and $\forall a \in F, OUT_{min}(a) = OUT_{max}(a) = 0,$
- where $InDegree(a)$ is $|\{e = (b, a) \mid e \in E\}|$ and $OutDegree(a)$ is $|\{e = (a, b) \mid e \in E\}|$.

Roughly speaking, an activity a can start as soon as at least $IN(a)$ of its predecessor activities have been completed. Two typical cases are: 1) if $IN(a) = InDegree(a)$, then a is an *and-join* activity, for it can be executed only after all its predecessors are completed and 2) if $IN(a) = 1$, it is called an *or-join* activity, for it can be executed as soon as one predecessor is completed. As is commonly assumed in the literature, we will limit ourselves to consider only *and-join* and *or-join* activities, besides a_0 : Indeed, by means of these two elementary types of nodes, it is also possible to simulate the behavior of any activity a such that

$$1 \leq IN(a) \leq InDegree(a).$$

Once finished, an activity a activates some (nondeterministically chosen) subset of its outgoing arcs with cardinality between $OUT_{min}(a)$ and $OUT_{max}(a)$. If $OUT_{max}(a) = OutDegree(a)$, then a is a *full fork* and if $OUT_{min}(a) = OUT_{max}(a)$ also, then a is a *deterministic fork*, for it activates all of its successor activities. Finally, if $OUT_{max}(a) = 1$, then a is an *exclusive fork* (also called *XOR-fork* in the literature), for it activates exactly one of their outgoing arcs.

For the sake of presentation, whenever it will be clear from the context, a workflow schema

1. <http://www.wfmc.org>.

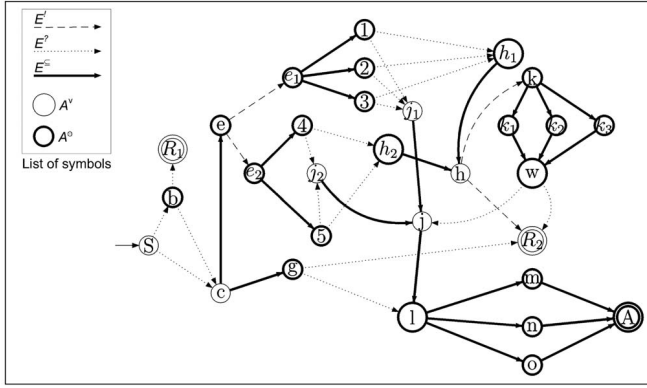


Fig. 1. An example workflow schema.

$$\mathcal{WS} = \langle A, E, a_0, F, \text{IN}, \text{OUT}_{\min}, \text{OUT}_{\max} \rangle$$

will also be denoted by $\langle A, E, a_0, F \rangle$ or, even simpler, by $\langle A, E \rangle$.

A workflow schema can be represented in a graphical way by means of a directed acyclic graph, where the nodes correspond to the activities in A and the edges correspond to the relation of precedence E (see Fig. 1). Moreover, in order to represent the functions IN , OUT_{\min} , and OUT_{\max} , if an activity is an *and-join* (respectively, *or-join*), we draw the corresponding node with bold (respectively, regular) circles; finally, the nodes corresponding to *exclusive fork* (respectively, *deterministic fork*) activities are such that their outgoing edges are marked with dotted (respectively, bold) lines, while all the other edges are represented by dashed lines.

Example 1. Fig. 1 shows a sketch of a workflow schema representing a sales ordering process. The process is as follows: A customer issues a request to purchase a given product; the enterprise checks both the availability of the required stock and the reliability of the client. Moreover, if the client is reliable but the products are partially stocked, a production will be planned. The final states can be the acceptance or the rejection of the order. Specifically, the initial task S corresponds to the “receive order” activity, the final tasks R_1 and R_2 are the rejecting of the order, and A is the acceptance. The activity e is the production that sends the request to some storehouse (either e_1 or e_2), which, in turn, forwards it to the respective repository (1, 2, 3 or 4, 5).

When at least one repository (no matter which one) has accepted the request, the task j_1 or j_2 proceeds to notify j . If there is no availability, the task h may send a request to the sales department (activity k) which forwards it to all wholesalers k_1, k_2, k_3 ; on the contrary, the user request will be rejected (task R_2).

Finally, the financial department (activity g) must assess if the reference is acceptable and if it is not, the order is rejected immediately (R_2); otherwise, the activation of the task l will lead to a success.

It is worth noting that, in this application, it could be crucial to characterize (with the help of the data mining techniques that we shall develop in the paper) the discriminant factors that will lead to an acceptance of the order requiring a planning of the production, in order to preventively accommodate the requests.

The formal semantics is specified by mapping the workflow schema into a transition system, where each execution consists of a sequence of states.

Definition 2. Let $\mathcal{WS} = \langle A, E, a_0, F \rangle$ be a workflow schema. Then, the state S of an execution is identified by a tuple $\langle \text{Marked}, \text{Ready}, \text{Executed} \rangle$, with $\text{Ready}, \text{Executed} \subseteq A$ and $\text{Marked} \subseteq E$.

Intuitively, the state of an execution is determined by the set (Executed) of activities which have been already executed by the set (Ready) of activities which have received the inputs they need and which are, hence, ready for being executed, and by the set (Marked) of edges corresponding to the outputs of executed activities which will eventually be inputs to other activities. An execution is modeled by means of a transition system over such states. Then, if after t transitions (short: step t), the state is $S_t = \langle \text{Marked}_t, \text{Ready}_t, \text{Executed}_t \rangle$ and the next state S_{t+1} is one of the outcomes of a nondeterministic transition function δ , defined next.

Definition 3. Let $\mathcal{WS} = \langle A, E, a_0, F \rangle$ be a workflow schema, and $S_t = \langle \text{Marked}_t, \text{Ready}_t, \text{Executed}_t \rangle$ be the state at the step t . Then, $\delta_{ws}(S_t)$ is the set of all states $\langle \text{Marked}_{t+1} \cup \delta \text{Marked}_{t+1}, \text{Ready}_{t+1}, \text{Executed}_t \cup \text{Ready}_t \rangle$, such that

1. $\delta \text{Marked}_{t+1}$ is a subset X of $\{(a, b) \mid a \in \text{Ready}_t, (a, b) \in E\}$ such that $\forall a \in \text{Ready}_t, \text{OUT}_{\min}(a) \leq |\{(a, b) \mid (a, b) \in X\}| \leq \text{OUT}_{\max}(a)$, i.e., each ready activity, say a , activates a number of outgoing arcs in the range defined by $\text{OUT}_{\min}(a)$ and $\text{OUT}_{\max}(a)$.
2. $\text{Ready}_{t+1} = \{a \mid a \in (A - (\text{Executed}_t \cup \text{Ready}_t)), |\{(b, a) \mid (b, a) \in \text{Marked}_t\}| \geq \text{IN}(a)\}$, i.e., an activity a becomes ready for execution as soon as at least $\text{IN}(a)$ of its predecessor activities are completed.

Now, we are in the position to formally define a workflow execution. An execution starts with the state $S_0 = \langle \emptyset, \{a_0\}, \emptyset \rangle$ and, at each step, it applies the transition function δ_{ws} until a final state is reached.

Definition 4. Let $\mathcal{WS} = \langle A, E, a_0, F \rangle$ be a workflow schema and δ_{ws} be a transition function. An execution e on a workflow schema $\mathcal{WS} = \langle A, E, a_0, F \rangle$ is a sequence of states $[S_0, \dots, S_k]$ such that

1. $S_0 = \langle \{\emptyset\}, \{a_0\}, \{\emptyset\} \rangle$ and
2. $S_{t+1} \in \delta(S_t)$ for each $0 \leq t \leq k$.

Moreover, if $\text{Executed}_k \cap F \neq \emptyset$ or $\text{Ready}_k \cup \delta \text{Marked}_k = \emptyset$, then e is said to be *terminating*; otherwise, it is said to be *partial*.

Given an instance $e = [S_0, \dots, S_k]$, the set Executed_k is also denoted by $\text{Executed}(e)$. Note that, in the above definition, a *terminating execution* e for which $\text{Executed}(e) \cap F = \emptyset$ corresponds to an abnormal execution which does not reach a final state. In this case, there are neither activities ready for being executed (i.e., $\text{Ready}_k = \emptyset$) nor outputs which may eventually activate other activities (i.e., $\delta \text{Marked}_k = \emptyset$); hence, e is said to be *unsuccessful*. Otherwise, e is said to be *successful*—observe that a successful execution may terminate with some ready activity that will be never executed, i.e., with $\text{Ready}_k \neq \emptyset$.

From now on, given a workflow schema \mathcal{WS} , the set of all the successful executions is denoted by S_{ws} , while the set of all the unsuccessful executions is denoted by U_{ws} .

Example 2. An example of execution over the workflow schema presented in Example 1 is reported in Fig. 2. The indexed columns represent the steps of the execution. Note that, at the fifth step, the financial department

step(t)	0	1	2	3	4	5	6	7
δMarked_t		(S, c)		$(c, e); (c, f); (c, g)$		$(e, e_1); (e, e_2); (g, R_2)$		$(e_1, 1); (e_1, 2); (e_1, 3); (e_2, 4); (e_2, 5)$
Ready_t	S		c		$e; f; g$		$e_1; e_2; R_2$	
Executed_t		S	S	$S; c$	$S; c$	$S; c; e; f; g$	$S; c; e; f; g$	$S; c; e; f; g; e_1; e_2; R_2$

Fig. 2. Example of execution over the workflow of Example 1

(activity g) has rejected the order (that is, not forwarded it to l), causing the ending of the workflow execution.

As suggested by the previous example, the choices made during an execution may cause a success or a failure. Moreover, checking whether the workflow has a sequence of choices leading to a success is an intractable problem. Specifically, we next show that it is complete for the class **NP** of problems that are solvable in polynomial time by nondeterministic Turing machines—see [28] for some background on computational complexity.

Proposition 1. *Let $\mathcal{WS} = \langle A, E, a_0, F \rangle$ be a workflow schema. Then, 1) deciding whether there exists an execution e that reaches a final state (i.e., $\text{Executed}(e) \cap F \neq \emptyset$) is **NP**-complete, but 2) the problem becomes **P**-complete if all nodes in A are full forks.*

Proof.

1. Membership in **NP** is trivial. For the hardness, recall that, given a Boolean formula Φ over variables X_1, \dots, X_m , the problem of deciding whether it is satisfiable is **NP**-complete [28]. Without loss of generality, assume Φ to be in conjunctive normal form. Then, we define a workflow schema $\mathcal{WS}(\Phi) = \langle A, E, a_0, \{Sat\} \rangle$, such that A consists of an initial activity a_0 of the activities X_i, TX_i, FX_i for each $0 \leq i \leq m$, the activities C_j for each distinct clause j of Φ , and a final state Sat . Moreover, we define $\text{IN}(Sat) = n$ (where n is the number of clauses contained in Φ) and $\text{IN}(a) = 1$ for any other activity $a \neq a_0$.

The set of precedences E is defined as follows:

- For each X_i , (X_i, TX_i) and (X_i, FX_i) are in E with constraints $\text{OUT}_{\min}(X_i) = \text{OUT}_{\max}(X_i) = 1$. Thus, each time the activity X_i is executed, it is required to make a choice between its possible successors; note that, in our encoding, TX_i means that X_i is **true**, while FX_i means that X_i is **false**. Finally, an arc (a_0, X_i) is in E and constraints $\text{OUT}_{\min}(a) = \text{OUT}_{\max}(a) = m$ are added.
- For each C_j , we have that (C_j, Sat) is in E with constraints $\text{OUT}_{\min}(Sat) = \text{OUT}_{\max}(Sat) = 1$. Moreover, we have that $(TX_i, C_j) \in E$ in the case that X_j appears in the clause C_j , while we have $(FX_i, C_j) \in E$ in the case that X_i appears negated in the clause C_j . Finally, for each node $a \in \{TX_i, FX_i\}$, $\text{OUT}_{\min}(a) = \text{OUT}_{\max}(a) = \text{OutDegree}(a)$.

Now, assume Φ is satisfiable. Then, it is possible to choose the successor of each X_i in a way that all the activities C_j can be executed. Hence, the activity Sat will eventually be reached. On the other side, if there exists a path leading to Sat , it can be easily mapped into a satisfying assignment for Φ .

2. Assume that, for each $a \in A$, $\text{OUT}_{\max}(a) = \text{OutDegree}(a)$. It is easy to see that, for the problem of deciding whether a given activity can be executed, we can assume without loss of generality that $\text{OUT}_{\min}(a) = \text{OutDegree}(a)$ too. Indeed, it is always convenient to activate all of the outgoing arcs in a : If an activity cannot be executed with the activation of all the outgoing arcs in a , then it cannot be executed in any other type of execution. Then, the problem can be solved in polynomial time by applying the function δ_{ws} , that it is actually a deterministic function. For the hardness, we consider the AND/OR GRAPH ACCESSIBILITY problem [29]: We are given an and/or graph $G = (V, E)$ (i.e., a directed graph such that each vertex is assigned either a \vee or \wedge label) and two vertices s and t ; the problem is to decide whether t can be reached from s . A vertex labeled by \vee can be reached if and only if at least one of its predecessors are reached, whereas a vertex labeled by \wedge can be reached if and only if all of its predecessors are reached. Notice that vertices without predecessors can be reached by default. We construct a workflow schema $\mathcal{WS}(G)$ by adding a starting activity connected to all the vertices without predecessors and s as well. For a vertex a labeled by \vee in G , we fix $\text{IN}(a) = 1$, whereas $\text{IN}(a) = \text{InDegree}(a)$ if a is labeled by \wedge in G . The only final activity is t . It is worth noticing that t is reachable from s in G if and only if $\mathcal{WS}(G)$ admits an execution reaching the final state. \square

3 PROBLEM DESCRIPTION AND COMPLEXITY RESULTS

In this section, we are interested in formalizing and analyzing the complexity of some interesting reasoning tasks whose usage can help the system administrator in predicting the workflow evolution. The analysis is carried out on the basis of the formal model of workflow schema and execution provided so far.

Let us first of all address the following problem: Assume that an execution has arrived at a given point and, before letting it proceed, the administrator wants to know whether it will lead to a successful termination or not. The problem can be formalized as follows.

Let \mathcal{WS} be a workflow schema and $e = [S_0, \dots, S_h]$ be a partial execution on \mathcal{WS} . A successful execution $e' \in \mathcal{S}_{ws}$ whose first $h+1$ steps are $[S_0, \dots, S_h]$ is said to be a *successful extension* of e and is denoted by $e \rightarrow e'$.

Definition 5 (Successful Termination Prediction—STP).

Let \mathcal{WS} be a workflow schema and e be a partial execution. Then, the STP problem for e is deciding whether there exists a successful extension of e .

We point out that the STP problem appears in [10] in the form of guaranteed executions in the more complex setting of transactional processes. We next show that STP is intractable.

Proposition 2. Let WS be a workflow schema and $e = [S_0, \dots, S_h]$ be an execution that is not terminating. Then, the STP problem for e is **NP-complete**.

Proof. Membership derives from the fact that we can guess an execution e' such that e' is successful and check (in polynomial time in the size of WS) whether its first $h+1$ steps are $[S_0, \dots, S_h]$. For the hardness, let $h=0$ and, without loss of generality, assume in the workflow schema the initial activity does not correspond to a final one. Thus, we can consider the problem of deciding whether there exists a successful execution e' with $[S_0] \rightarrow e'$. The hardness follows from Proposition 1. \square

The above discussion sheds some light in the intrinsic difficulty of solving such problems “statically.” Reasoning about the structure seems not to be a valuable approach; hence, we are motivated in using data mining techniques that can be directly applied to a set of instances collected in the log of the workflow system. Indeed, one could be interested in a more pragmatic version of the STP problem: *Given the history of past executions, does the current execution have a chance to eventually succeed?* We formalize the problem next.

Definition 6 (Frequent Successful Termination Prediction—FSTP). Let WS be a workflow schema, $Se = \{e_1, \dots, e_n\}$ be a set of successful executions on WS , each one equipped with a frequency $f_i = f_i(e_i) \in \mathbb{N}$, $minFreq$ be a natural number, and e be a nonterminating execution on WS . Then, the problem FSTP for e with regard to Se and $minFreq$ is deciding whether $\sum_{\{i|e_i \in Se, e \rightarrow e_i\}} f_i \geq minFreq$, i.e., whether the number of successful extensions of e in Se is greater than or equal to $minFreq$.

As a matter of fact, the STP problem is equivalent to an instance of the FSTP problem.

Proposition 3. Let a workflow schema WS and a nonterminating execution e be given in input. Then, the STP problem is equivalent to the FSTP problem for e with regard to S_{ws} , where $minFreq = 1$ and $f_i(e_i) = 1$ for each execution $e_i \in S_{ws}$.

Proof. By definition, under the assumptions of the statement, the problem FSTP corresponds to check whether $|\{e_i | e_i \in S_{ws}, e \rightarrow e_i\}| \geq 1$. This happens if and only if there exists $e_i \in S_{ws}$ such that $e \rightarrow e_i$. \square

The complexity of the FSTP problem mainly depends on the number of executions in Se . If this number is low (e.g., polynomially bounded by the size of WS), then the problem can be effectively solved, as the following proposition shows. Nevertheless, when the size of Se grows, one cannot expect to reduce the complexity by finding some succinct representation of Se : Also, in this case, a time exponential in the size of WS cannot be avoided unless $P = NP$.

Proposition 4. Let WS be a workflow schema and $e = [S_0, \dots, S_h]$ be an execution that is not terminating. Then, given a set $Se = \{e_1, \dots, e_n\}$ of terminating executions on WS , each one equipped with a frequency $f_i = f_i(e_i) \in \mathbb{N}$ and a natural number $minFreq$,

- the FSTP problem for e with regard to Se and $minFreq$ can be solved in time polynomial in the size of WS and Se , but

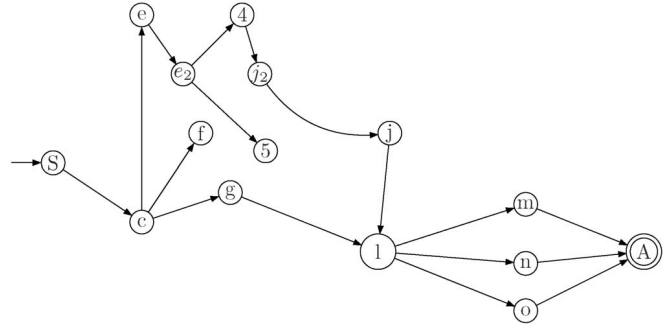


Fig. 3. An instance of the workflow schema of Fig. 1.

- the succinct FSTP problem, in which Se is represented by a data structure with size polynomially bound in the size of WS , is **NP-complete**.

Proof.

1. Observe that, for each e_i , checking whether $e \rightarrow e_i$ can be done in polynomial time in the size of WS . Hence, a naive polynomial algorithm (in the size of WS and Se) consists in summing the frequency associated to each execution $e_i \in Se$ with $e \rightarrow e_i$ and, hence, checking whether the corresponding sum is greater than $minFreq$.
2. Membership is trivial. For the hardness, observe that the FSTP problem in the statement of Proposition 3 can be obviously formulated in a succinct way: Indeed, Se needs not to be explicitly stored. But, the succinct problem is nothing but the STP problem, which is **NP-complete** by Proposition 2. \square

An appealing way for solving the FSTP problem is to use specialized data mining techniques for graphs. To this end, we first need to characterize workflow executions in terms of connected subgraphs of the workflow schema.

Definition 7. Let $WS = \langle A, E, a_0, F \rangle$ be a workflow schema and $e = [S_0, \dots, S_k]$ be an execution. Then, the instance associated to e is the graph $I_e = \langle A_e, E_e, a_0, F_e \rangle$, where $A_e = \cup_{t=1,k} Executed_t$, $E_e = \{(a, b) | (a, b) \in \cup_{t=1,k} Marked_t, b \in A_e\}$, and $F_e = A_e \cap F$. In case e is a successful execution, I_e is said to be a successful instance.

An instance for the workflow schema presented in Example 1 is shown in Fig. 3.

In the following, given a workflow schema WS , we denote by 2^{WS} the family of all the subgraphs of the graph $\langle A, E \rangle$ and by $\mathcal{I}(WS)$ the set of all instances.

Observe that, while deciding whether a subgraph is an instance that is polynomial, instead deciding whether there exists a successful instance is not tractable.

Proposition 5. Let $WS = \langle A, E, a_0, F \rangle$ be a workflow schema. Then,

- given a subgraph I of WS , deciding whether I is an instance of WS can be done in polynomial time in the size of E and
- deciding whether WS admits a successful instance is **NP-complete**.

Proof.

1. We construct the sequence of states corresponding to I by traversing the subgraph I starting from the initial node and by applying in a constructive way the function δ using all arcs in I as marked. Clearly, the algorithm is polynomial in the size of E .
2. Membership is trivial; for the hardness, observe that there exists a successful instance, say I_e , if and only if there exists a successful execution, e , to which I_e is associated. The hardness follows from Proposition 1. \square

We now introduce the notion of pattern that will be crucial in the process of data mining. To provide a more uniform notation, given a graph p (e.g., a workflow schema or a pattern) and a node a of p , we denote by $InDegree_p(a)$ (respectively, $OutDegree_p(a)$) the number of ingoing (respectively, outgoing) edges of a .

Definition 8. Let \mathcal{WS} be a workflow schema and \mathcal{F} be a multiset of instances. Then, a graph $p = \langle A_p, E_p \rangle \in 2^{\mathcal{WS}}$ is an \mathcal{F} -pattern (cf. $\mathcal{F} \models p$) if there exists $I = \langle A_I, E_I \rangle \in \mathcal{F}$ such that $A_p \subseteq A_I$ and p is the subgraph of I induced by the nodes in A_p . In the case $\mathcal{F} = \mathcal{I}(\mathcal{WS})$, the subgraph is simply said to be a pattern. Moreover, if A_p contains some final activity in \mathcal{WS} , then p is said to be successful.

Roughly speaking, an \mathcal{F} -pattern is a subgraph of a workflow instance in \mathcal{F} . Thus, we are using the notion of pattern with the meaning which is being adopted by the data mining community in several other applicative domains. For instance, patterns are subtrees in the mining of frequent trees (see, e.g., [2]), subsequences in the mining of sequences (see, e.g., [1]), and so on. Hence, Definition 8 is inserted into a data mining context and is not related at all with the notion of pattern used in software engineering contexts (and recently by van der Aalst, et al. [30] for supporting workflow modeling²).

Let us now consider the following problem of data mining on graphs that consists in discovering patterns which frequently arise.

Definition 9 (Frequent pattern mining—FPM). Let \mathcal{WS} be a workflow schema, \mathcal{F} a multiset of instances, and $minSupp$ a real number with $0 \leq minSupp \leq 1$. Then, the problem FPM for \mathcal{F} consists in finding all the frequent \mathcal{F} -patterns, i.e., all the \mathcal{F} -patterns for which $supp(p) \geq minSupp$, where the support $supp(p)$ is defined as $|\{I \in \mathcal{F} \mid I \models p\}|/|\mathcal{F}|$.

Frequent patterns can be used for heuristically solving the problem FSTP, that is, for deciding whether a sequence of states will very likely (e.g., with a reasonable support) lead to a successful (or unsuccessful) termination. In fact, given a partial execution e and a support $minSupp$, e will very likely lead to a successful end if there exists at least one successful frequent \mathcal{F} -pattern containing all the activities executed in e . Then, in order to make our approach effective, we will show in the following section some techniques for the efficient computation of frequent patterns of executions.

2. See also <http://tmitwww.tn.tue.nl/research/patterns>.

4 MINING CONNECTED FREQUENT PATTERNS

In this section, we present two algorithms for mining connected frequent patterns (i.e., subgraphs) in workflow instances. Let us assume that a workflow schema $\mathcal{WS} = \langle A, E, a_0, F \rangle$ and a multiset of instances $\mathcal{F} = \{I_1, \dots, I_n\}$ are given. Then, a naive algorithm for mining frequent patterns can generate directly the subgraphs and test in polynomial time whether they are instances of \mathcal{WS} . Our approach is based on the idea of reducing the number of patterns to generate by only considering \mathcal{F} -patterns that are not only connected but also *deterministically closed*. This restriction is formalized next.

Definition 10. Given a graph $p = \langle A_p, E_p \rangle \in 2^{\mathcal{WS}}$, the *deterministic closure* of p (cf. $ws\text{-closure}(p)$) is inductively defined as the graph $p' = \langle A_{p'}, E_{p'} \rangle$ such that:

1. $A_p \subseteq A_{p'}$ and $E_p \subseteq E_{p'}$ (basis of induction),
2. $a \in A_{p'}$ is an and-join implies that for each $(b, a) \in E$, $(b, a) \in E_{p'}$ and $b \in A_{p'}$,
3. $a \in A_{p'}$ is a deterministic fork implies that for each $(a, b) \in E$ with b or-join,³ $(a, b) \in E_{p'}$ and $b \in A_{p'}$. Moreover, a graph p such that $p = ws\text{-closure}(p)$ is said to be *ws-closed*.

Intuitively, the above definition provides a way for extending a subgraph p by including all the activities that are enforced to be executed with some activity in A_p by means of the constraints issued over \mathcal{WS} . And, in fact, the definition can be used to introduce a notion of pattern which only depends on the structure of the workflow schema, rather than on the instances \mathcal{F} or $\mathcal{I}(\mathcal{WS})$. The need of this weaker notion will be clear in a while.

Definition 11. A weak pattern, or simply *w-pattern*, is a ws-closed graph $p \in 2^{\mathcal{WS}}$, such that for each node a , $|\{(a, b) \mid (a, b) \in E_p\}| \leq OUT_{max}(a)$.

Example 3. Consider the workflow graph of Fig. 1 and the subgraphs in Fig. 4. Then, p_1 is not a *w-pattern* since $ws\text{-closure}(p_1) = p_2 \neq p_1$ and, hence, condition 2 of Definition 10 is not satisfied. Notice that p_2 is instead a *w-pattern* since $ws\text{-closure}(p_1) = p_2$. Also, p_3 is not a *w-pattern* since condition 2 of Definition 10 is not satisfied (indeed, $ws\text{-closure}(p_3) = p_4 \neq p_3$). Again, p_4 is a *w-pattern*, as $ws\text{-closure}(p_4) = p_4$.

The following proposition characterizes the complexity of recognition for the three notions of pattern; in particular, it states that testing whether a graph is a *w-pattern* is in **L** [28], i.e., it can be efficiently solved by a deterministic logarithmic-space bounded Turing machine. This efficiency is the result of the deterministic closure property and of the fact that *w-patterns* are defined over the schema, rather than on the instances.

Proposition 6. Let $p \in 2^{\mathcal{WS}}$. Then,

1. Deciding whether p is a pattern is **NP-complete**.
2. Given a multiset \mathcal{F} of instances, deciding whether p is an \mathcal{F} -pattern can be done in polynomial time in the size of \mathcal{F} .

3. Notice that relaxing the condition for b to be an *or-join* might lead to closures that cannot be traced by any execution. In fact, if b is an *and-join* synchronizing two mutually exclusive activities a and a' (e.g., that are activated by some *XOR-fork*), then a will never occur in the same execution with b .

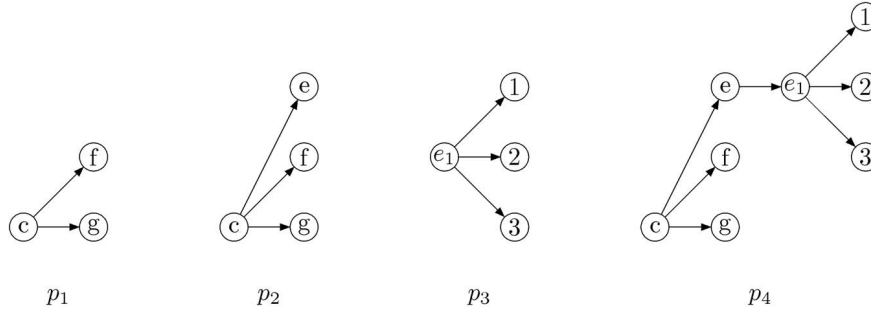


Fig. 4. Example subgraphs of the schema of Fig. 1. p_2 and p_4 are w -patterns.

3. Deciding whether p is an \mathcal{F} -pattern is **NP-complete** if \mathcal{F} is succinct (i.e., if it can be represented by a data structure whose size is polynomially bounded by the size of \mathcal{WS}).
4. Deciding whether p is a w -pattern is in **L**.

Proof.

1. The problem is in **NP** as we can guess a subgraph I by choosing for each node a the arcs to be activated so that $\text{OUT}_{\min}(a) \leq |\{(a,b) | (a,b) \in E_p\}| \leq \text{OUT}_{\max}(a)$. Then, from Proposition 5, we can check in polynomial time that I is an instance; finally, deciding whether p is a subgraph of I can be done in polynomial time.

The hardness follows from Proposition 1; indeed, we can assume p to be formed by a single activity, actually a final one (without loss of generality, we can assume that \mathcal{WS} has only one final activity, indeed we can add a new activity f to which all the final ones can be connected). Thus, p is a pattern if and only if there exists a successful execution.

2. By definition of \mathcal{F} -pattern, we can simply test if p is a subgraph of any instance.
3. Membership is trivial, as we can check in polynomial time whether, for each instance I , $\{I\} \models p$ and $I \in \mathcal{F}$. For the hardness, observe that deciding whether p is a pattern corresponds to checking whether p is a \mathcal{F} -pattern with $\mathcal{F} = \mathcal{I}(\mathcal{WS})$.
4. The proof is given by defining a Turing machine that, given a workflow schema and a graph p encoded into the input tapes, can decide in deterministic logarithmic space whether p is a w -pattern. In fact, both \mathcal{WS} and p can be encoded by fixing an arbitrary order on the activities. In order to verify properties 1, 2, and 3 of Definition 10, we simply need to access each arc of p and \mathcal{WS} and exploit two counters. Clearly, encoding such counters requires logarithmic space. \square

It turns out that the notion of weak pattern is the most appropriate from the computational point of view. Moreover, working with w -patterns is not an actual limitation, since the closure of each frequent \mathcal{F} -pattern is, in turn, a frequent w -pattern as well. Rather, it is a compact and efficient way for the mining of frequent patterns, as shown below.

Proposition 7. Let p be a frequent \mathcal{F} -pattern. Then, 1) $\text{ws-closure}(p)$ is both a weak pattern and a frequent \mathcal{F} -pattern and 2) each weak pattern $p' \subseteq p$ is a frequent \mathcal{F} -pattern.

Proof. In order to prove 1), we observe that, for each $I \in \mathcal{F}$ such that $\{I\} \models p$, property $\{I\} \models \text{ws-closure}(p)$ holds. Indeed, if p is not a weak pattern, then, according to Definition 10, there exists $a \in A_p$ such that one of the following cases occur:

- a is an *and-join* and there exists an edge $(b,a) \notin E$;
- a is a *deterministic fork* and there exists an edge $(a,c) \notin E_p$ with c *or-join*.

By Definitions 3 and 7, each instance $I \in \mathcal{F}$ containing a , must contain b, c and $(b,a), (a,c)$ as well. As a consequence, $\text{ws-closure}(p)$ is frequent as well.

In order to prove 2), it suffices to see that, if there exists an unfrequent w -pattern $p' \subseteq p$, then it should contain at least either an unfrequent node a or an unfrequent edge (a,b) . But, this is a contradiction since both a and (a,b) belong to p as well. \square

However, a weak pattern is not necessarily an \mathcal{F} -pattern nor even a pattern. As shown in the next sections, we shall use weak patterns in our mining algorithms to optimize the search space, but we eventually check whether they are frequent \mathcal{F} -patterns.

4.1 Levelwise Search Algorithm

The first algorithm we propose for mining frequent connected \mathcal{F} -patterns uses a levelwise theory. Roughly speaking, we incrementally construct frequent weak patterns by starting from frequent “elementary” weak patterns (defined below) and by extending each frequent weak pattern using two basic operations: adding a frequent arc and merging with another frequent elementary weak pattern. As we shall show, the correctness follows from the results of Proposition 7 and from the observation that the space of all connected weak patterns constitutes a lower semilattice with a precedence relation \prec , defined next.

The elementary weak patterns from which we start the construction of frequent patterns are obtained as the ws-closures of the single nodes.

Definition 12. Let $\mathcal{WS} = \langle A, E \rangle$ be a workflow schema. Then, for each $a \in A$, the graph $\text{ws-closure}(\{a\}, \{\})$ is called an elementary weak pattern (cf. *ew-pattern*).

Observe that the empty graph, denoted by \perp , is an elementary weak pattern. The set of all *ew*-patterns is denoted by EW . Moreover, let p be a weak pattern, then EW_p denotes the set of the elementary weak patterns contained in p . Note that, given an *ew*-pattern e , EW_e is not necessarily a singleton, for it may contain other *ew*-patterns.

Given a set $E' \subseteq \text{EW}$, $\text{Compl}(E') = \text{EW} - \bigcup_{e \in E'} \text{EW}_e$ contains all elementary patterns which are neither in E' nor contained in some element of E' .

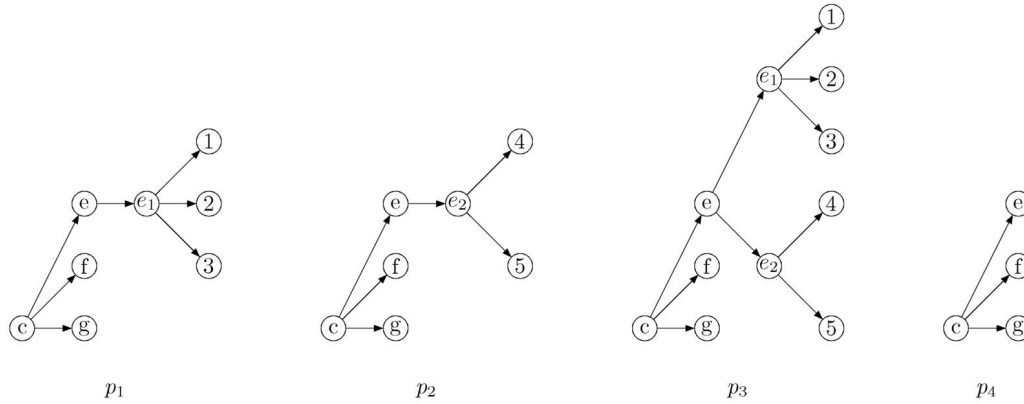


Fig. 5. Example elementary patterns and precedence relations.

We now introduce a precedence relation \prec among connected weak patterns. First of all, let us denote by E^\subseteq the subset of arcs in \mathcal{WS} whose source is not a deterministic fork, i.e., $E^\subseteq = \{(a, b) \in E \mid \text{OUT}_{\min}(a) < \text{OutDegree}_p(a)\}$.

Definition 13. Given two connected w -patterns, say $p = \langle A_p, E_p \rangle$ and $p' = \langle A_{p'}, E_{p'} \rangle$, $p \prec p'$ if and only if:

1. $A_p = A_{p'}$ and $E_{p'} = E_p \cup \{(a, b)\}$, where $(a, b) \in E^\subseteq - E_p$ and $\text{OUT}_{\max}(a) > \text{OutDegree}_p(a)$ (i.e., p' can be obtained from p by adding an arc), or
2. there exists $p'' \in \text{Compl}(\text{EW}_p)$ such that $p' = p \cup p'' \cup X$, where X is either empty if p and p'' are connected or contains exactly an edge in E^\subseteq with endpoints in p and p'' (i.e., p' is obtained from p by adding an elementary weak pattern and possibly a connecting arc).

Note that $\perp \prec e$ for each $e \in \text{EW}$.

Example 4. Consider again the workflow of Fig. 1 and the subgraphs in Fig. 5. The subgraphs p_1 , p_2 , and p_4 are elementary patterns: $p_1 = \text{ws-closure}(\langle \{e_1\}, \emptyset \rangle)$, $p_2 = \text{ws-closure}(\langle \{e_2\}, \emptyset \rangle)$, and $p_4 = \text{ws-closure}(\langle \{c\}, \emptyset \rangle)$. p_3 is not an elementary pattern as no node can generate it. Notice that $p_1 \prec p_3$ and $p_2 \prec p_3$. Finally, p_4 is contained in both p_1 and p_2 (and, hence, $p_4 \prec p_1$ and $p_4 \prec p_2$).

The following result states that all the connected weak patterns of a given workflow schema can be constructed by means of a chain over the \prec relation.

Lemma 1. Let $p = \langle A_p, E_p \rangle$ be a connected w -pattern. Then, there exists a chain of connected w -patterns, such that $\perp \prec p_1 \prec \dots \prec p_n = p$.

Proof. We prove this by induction on the size of p , $|p| = |A_p| + |E_p|$. The base case, i.e., $p \in \perp$, is trivial. For the case $p \notin \text{EW}$, assume that for each weak pattern p' such that $|p'| \leq |p|$, there exists a chain $\perp \prec q_1 \prec \dots \prec q_m = p'$.

Two situations may occur:

1. $\exists (a, b) \in E_p \cap E^\subseteq$, such that (a, b) does not belong to any elementary pattern contained in p and the graph p' obtained from p by deleting such an arc ($p' = \langle A_p, E_p - \{(a, b)\} \rangle$) is connected. In such a case, p' is a weak pattern with $p' \prec p$. Hence, by induction, $\perp \prec q_1 \prec \dots \prec q_m \prec p$. The theorem follows for $n = m$ and $p_1 = q_1, \dots, p_n = q_n$.
2. For each $(a, b) \in E_p \cap E^\subseteq$, such that (a, b) does not belong to any elementary pattern contained in p ,

the graph $p' = \langle A_p, E_p - \{(a, b)\} \rangle$ is not connected. Two subcases can be further devised:

- a. There exists an elementary weak pattern $e \in \text{EW}_p$, which is connected to the graph $p - e$ by means of exactly one arc in E^\subseteq ; that is, $e \in \text{Compl}(\text{EW}_{p-e})$ and, hence, $(p - e) \prec p$ and the theorem follows by induction.
- b. Elementary patterns are not connected by means of arcs in E^\subseteq . In this case, let ep_0, ep_1, \dots, ep_m be the elementary patterns contained in p and $q = (p - ep_0) \cup ep_1 \cup \dots \cup ep_m$ be the weak pattern obtained from p by deleting edges and nodes in ep_0 which do not occur in any other ep_i with $0 \leq i \leq m$. By construction, $ep_0 \in \text{Compl}(q)$ and, hence, $q \prec p$. As in the other case, since $|q| \leq |p|$, by induction there exists a chain of weak patterns $\perp \prec q_1 \prec \dots \prec q_m = q$. \square

It turns out that the space of all connected weak patterns is a lower semilattice with regard to the precedence relation \prec . And, in fact, the algorithm $w\text{-find}$, reported in Fig. 6, exploits an a-priori-like exploration of this lower semilattice. Specifically, at each stage, the computation of L_{k+1} (steps 5-14) is carried out by extending any pattern p generated at the previous stage ($p \in L_k$) in two ways: 1) by adding frequent edges in E^\subseteq (addFrequentArc function) and 2) by adding an elementary weak pattern ($\text{addEWFrequentPattern}$ function). Each pattern p' , generated by the functions above, is an admissible subgraph of \mathcal{WS} (cf. $\mathcal{WS} \models p'$), i.e., for each $a \in A_{p'}$, $\text{OutDegree}_{p'}(a) \leq \text{OUT}_{\max}(a)$. The properties of the $w\text{-find}$ algorithm are reported below.

Lemma 2. In the $w\text{-find}$ algorithm, the following propositions hold:

1. addFrequentArc adds to U connected patterns, which are not necessarily \mathcal{F} -patterns.
2. $\text{addFrequentEWPattern}$ adds to U connected w -patterns (not necessarily patterns).
3. For each k , L_k contains only frequent connected \mathcal{F} -patterns.

Proof. We shall prove the above statements by induction on k . The proof is structured as follows: First, observe that L_0 contains a set of frequent connected \mathcal{F} -patterns. Indeed, by definition, each elementary weak pattern is

Input: A workflow Graph \mathcal{WS} , a set $\mathcal{F} = \{I_1, \dots, I_N\}$ of instances of \mathcal{WS} .
Output: A set of frequent \mathcal{F} -patterns.
Method: Perform the following steps:

```

1   $L_0 := \{e | e \in EW, e \text{ is frequent w.r.t. } \mathcal{F}\};$ 
2   $k := 0, R := L_0;$ 
3   $FrequentArcs := \{(a, b) | (a, b) \in E^\subseteq, \langle \{a, b\}, \{(a, b)\} \rangle \text{ is frequent w.r.t. } \mathcal{F}\};$ 
4   $E_f^\subseteq := E^\subseteq \cap FrequentArcs;$ 
5  repeat
6     $U := \emptyset;$ 
7    forall  $p \in L_k$  do begin
8       $U := U \cup addFrequentArc(p);$  //see (a) in Lemma 2
9      forall  $e \in Compl(EW_p) \cap L_0$  do
10        $U := U \cup addFrequentEWPattern(p, e);$  //see (b) in Lemma 2
11    end
12     $L_{k+1} := \{p | p \in U, p \text{ is frequent w.r.t. } \mathcal{F}\};$  //see (c) in Lemma 2
13     $R := R \cup L_{k+1};$ 
14  until  $L_{k+1} = \emptyset;$ 
15  return  $R;$ 

```

Function $addFrequentEWPattern(p = \langle A_p, E_p \rangle, e = \langle A_e, E_e \rangle)$: **w-pattern**;
 $p' := \langle A_p \cup A_e, E_p \cup E_e \rangle;$
if p' **is connected**, **then return** p' **else return** $addFrequentConnection(p', p, e);$

Function $addFrequentConnection(p' = \langle A_{p'}, E_{p'} \rangle, p = \langle A_p, E_p \rangle, e = \langle A_e, E_e \rangle)$: **w-pattern**;
 $S := \emptyset$
forall frequent $(a, b) \in E_f^\subseteq - E_p$ s.t. $(a \in A_p, b \in A_e) \vee (a \in A_e, b \in A_p)$ **do begin**
 $q := \langle A_{p'}, E_{p'} \cup (a, b) \rangle;$
if $\mathcal{WS} \models q$ **then** $S := S \cup \{q\};$
end
return S

Function $addFrequentArc(p = \langle A_p, E_p \rangle)$: **pattern**;
 $S := \emptyset$
forall frequent $(a, b) \in E_f^\subseteq - E_p$ s.t. $a \in A_p, b \in A_p$ **do begin**
 $p' := \langle A_p, E_p \cup (a, b) \rangle;$
if $\mathcal{WS} \models p'$ **then** $S := S \cup \{p'\};$
end
return S

Fig. 6. Algorithm $w\text{-find}(\mathcal{F}, \mathcal{WS})$.

connected. Next, assuming that for a given $k \geq 0$, L_k contains only a set of frequent connected \mathcal{F} -patterns, observe that:

- Statement 1 holds. Indeed, since the input graph p is a connected \mathcal{F} -pattern, each graph p' obtained from p by adding a frequent arc (a, b) is connected as well. Notice now that, if $\mathcal{WS} \models p'$, then for each instance $\{I\} \models p$, the graph $I' = I \cup \{(a, b)\}$ is an admissible instance, i.e., $I' \in \mathcal{I}(\mathcal{WS})$. Indeed, for each execution e_I associated to I , an execution $e_{I'}$ can be obtained by adding (a, b) to $\delta Marked_{t+1}$ whenever $a \in Ready_t$. Moreover, $\{I'\} \models p'$, thus entailing that p' is a pattern. Finally, notice that, in principle, \mathcal{F} may contain no I' satisfying the above condition.
- Statement 2 holds. Indeed, notice that $addEW\text{-}FrequentPattern$ returns any admissible connected subgraph p obtained from the union of a \mathcal{F} -pattern p' with an elementary pattern p'' . If p is not a w -pattern, then either there exists an *and-join* $a \in A_p$ and $(b, a) \in E$ such that $(b, a) \notin E_p$, or there exists a *deterministic fork* $a \in A_p$ and $(a, b) \in E$ with b *or-join*, such that $(a, b) \notin E_p$. In both cases, there exists a node a and an edge e (containing a) such that $e \notin E_p$. But, this cannot happen, since a is contained either in p' or in p'' : Indeed, since p' and p'' are w -patterns, e is contained in any of

them and, consequently, in E_p . Finally, notice that p'' is not necessarily a pattern and, consequently, p is not necessarily a pattern as well.

- Statement 3 holds. Indeed, it follows from statements 1 and 2 that the set of candidate graphs U contains either connected patterns or connected \mathcal{F} -patterns. The consequence is trivial by noticing that step 12 of the algorithm adds to L_{k+1} the only patterns in U which are frequent with regard to \mathcal{F} . \square

We next show that all the weak patterns are actually computed by the algorithm.

Proposition 8 (Correctness). *The algorithm of Fig. 6 terminates and computes all and only the frequent connected weak patterns.*

Proof. The algorithm $w\text{-find}$ computes all the elements in the lower semilattice induced by the operator \prec over w -patterns. The correctness follows from Lemma 1, stating that any weak pattern is represented by a chain in this lower semilattice, and by the observation that we also prune the chains that will lead to unfrequent patterns. The latter is done by replacing the function $addEWPattern$ in the definition of the relation \prec with $addFrequentEWPattern$. \square

As a conclusion of the presentation of $w\text{-find}$, we again remark that focusing on weak patterns is an efficient way for

computing frequent patterns. In fact, Proposition 7 states that 1) for each frequent \mathcal{F} -pattern p' , there exists a frequent weak pattern p (hence, computed by w -find) containing p' and 2) any subgraph of any frequent weak pattern (again, computed by w -find) is a frequent \mathcal{F} -pattern as well.

4.2 Mining by Connecting Components

The algorithm w -find, proposed in Fig. 6, performs a smart levelwise exploration of the lower semilattice analyzed in Lemma 1. However, a different strategy can be exploited by observing that, in general, any connected pattern can be obtained by either composing two connected subgraphs or by extending a subgraph by means of an edge.

Lemma 3. *Let $p \in (2^{\mathcal{WS}} - EW)$ be a connected \mathcal{F} -pattern. Then, there exist two \mathcal{F} -patterns p_1 and p_2 (not necessarily distinct) such that $p = p_1 \cup p_2 \cup X$, where X can be either the empty set or the graph $\langle \{a, b\}, \{(a, b)\} \rangle$ with $a \in p_1$ and $b \in p_2$.*

Proof. Let p be a connected pattern not in EW . Then, due to Proposition 7, $q = ws\text{-closure}(p)$ is a weak pattern. Due to Lemma 1, there exists a chain of connected w -patterns such that $q_0 = \perp < q_1 < \dots < q_{n-1} < q_n = q$. Moreover, each q_{i+1} can be derived from q_i by either adding an edge in p_i or by connecting an elementary weak pattern to q_i . By denoting with Δq_i the graph that we compose with q_i , in order to derive q_{i+1} , we derive the following relationship: $\Delta q_0 < \Delta q_0 \cup \Delta q_1 < \dots < \bigcup_{i=0, n-1} \Delta q_i = q$. Without loss of generality, we can assume that there exists $0 \leq j \leq n$ such that $p'_1 = \bigcup_{i=0, j} \Delta q_i$ and $p'_2 = \bigcup_{i=j+1, n-1} \Delta q_i$ are connected (it trivially holds for $j = n - 1$). Then, by the definition of the relation of precedence $<$, we have that p'_1 and p'_2 are either connected or can be connected by means of an edge. The result follows by letting $p_1 = p'_1 \cap p$ and $p_2 = p'_2 \cap p$. \square

The above lemma states that candidates can be generated by iteratively connecting components. In fact, we can generate a candidate at the n th level of the lattice by merging two components at the j th and $(n-j)$ th level, respectively. It is clear that in the worst case, for $j = n-1$, we degenerate to the levelwise search described in the previous section; nonetheless, in the best case, this approach converges in exponentially fewer iterations. Obviously, we also need an additional effort for identifying the components that can be merged. Roughly speaking, these components must be such that their boundaries can match, where the boundary of a graph in $2^{\mathcal{WS}}$ is the set of nodes that (according to a workflow schema \mathcal{WS}) admit either an input or an output edge.

In order to formalize the above intuitions, given a graph $p = \langle A_p, E_p \rangle \in 2^{\mathcal{WS}}$, we denote by $INBORDER(p) = \{a \in A_p \mid InDegree_p(a) < InDegree(a)\}$ the set of all the nodes in p which admit a further incoming edge and by $OUTBORDER(p) = \{a \in A_p \mid OutDegree_p(a) \leq OUT_{max}(a)\}$ the set of all nodes in p which admit an outgoing edge. The sets $INBORDER(p)$ and $OUTBORDER(p)$ represent the input and output boundaries of p , i.e., the set of nodes inside p , which can reach (respectively, can be reached by) other nodes outside p . Notice that, by construction, the input boundary of a w -pattern cannot contain *and-join* activities. Similarly, the output boundary of a w -pattern cannot contain *deterministic forks*.

The boundaries can be exploited to connect components. Since an arc connects the boundaries of two components, it suffices to concentrate on frequent arcs and iteratively

generate new candidates by merging the frequent components whose boundaries are connected by means of those arcs.

Based on this idea, we have developed another algorithm (c -find), whose details are reported in Fig. 7. It starts by computing frequent elementary patterns (step 1). Then, the core of the algorithm is a main loop (steps 3-24), in which the following operations are performed. For each node $a \in \mathcal{WS}$, the set $INF(a)$ (respectively, $OUTF(a)$) of \mathcal{F} -patterns containing a in the input (respectively, output) boundary is computed (steps 5-6). In steps 8 and 9, the variables FA and FP are used to store frequent arcs that may connect patterns and candidates that may be generated by composing “compatible” patterns.

Then, boundaries are recomputed for the new candidate \mathcal{F} -patterns (steps 11-21) and frequent \mathcal{F} -patterns are detected by computing the frequency of each candidate (step 22). Notice that boundaries for candidate \mathcal{F} -patterns can be incrementally computed by extending the boundaries of the connected components and that new candidates can be generated also by merging \mathcal{F} -patterns sharing some nodes. The algorithm terminates when no further candidates can be found, i.e., when the computed patterns have empty input-output boundaries.

Theorem 1 (Correctness). *The c -find algorithm terminates and computes all and only the frequent connected weak patterns.*

Proof. Correctness trivially holds by step 22 of the algorithm: Indeed, no pattern is included in R unless it is not an \mathcal{F} -pattern. As for completeness, let p be an \mathcal{F} -pattern. We prove by induction on $|p|$ that $p \in R$. The case $p \in EW$ statement trivially holds as a consequence of step 1 of the algorithm. Let us consider the case $|p| \geq 1$. By Lemma 3, there exist p_1, p_2 such that $p = p_1 \cup p_2 \cup X$, where X can be the empty set or an edge connecting p_1 and p_2 . By induction, both p_1 and p_2 are in R . Let us assume, without loss of generality, that p_1 is added to R at iteration k_1 and that p_2 is added to R at iteration $k_2 \geq k_1$. Two situations may occur:

1. $p_1 \cap p_2 \neq \emptyset$ and $p = p_1 \cup p_2$. In such a case, p is added to PF at the iteration $k_2 + 1$ and, consequently, it is added to R .
2. $p = p_1 \cup p_2 \cup \{(a, b)\}$. Again, without loss of generality, we can assume that $a \in OUTBORDER(p_1)$ and $b \in INBORDER(p_2)$. In such a case, by step 8 of the algorithm, $(a, b) \in FA$ at iteration $k_2 + 1$. By steps 11 and 12 of the algorithm, $p \in PF$ at iteration $k_2 + 1$ and, hence, $p \in R$.

Observe finally that each candidate pattern p is considered at most k times, where k is the number of connected patterns contained in p . Since the number of candidate patterns is finite, the algorithm must terminate. \square

In comparing the performance of the c -find algorithm with the w -find algorithm proposed in the previous section, it is interesting to notice that the c -find algorithm can generate more candidates than w -find, but in general reaches convergence more quickly (number of iterations).

Proposition 9. *Let C be the set of candidate patterns generated by c -find and let N_c be the steps required for its execution. Let \mathcal{W} be the set of candidate patterns generated by w -find and let N_w be the steps required for its execution. Then, $\mathcal{W} \subseteq C$ and $N_c \leq N_w$.*

Input: A workflow Graph $\mathcal{WS} = (A, E)$, a set $\mathcal{F} = \{I_1, \dots, I_N\}$ of instances of \mathcal{WS} .
Output: A set of frequent \mathcal{F} -patterns.
Method: Perform the following steps:

```

1   $R := \{e \mid e \in EW, e \text{ is frequent w.r.t. } \mathcal{F}\}; \Delta R := R;$ 
2  forall  $(a, b) \in E$  do  $connected\_by(a, b) = \emptyset;$ 
3  repeat
4    forall  $a \in A$  do begin
5       $INF(a) := \{p \in R \mid a \in INBORDER(p)\}; INFP(a) = \emptyset;$ 
6       $OUTF(a) := \{p \in R \mid a \in OUTBORDER(p)\}; OUTFP(a) = \emptyset;$ 
7    end
8     $FA := \{(a, b) \mid (a, b) \text{ is frequent w.r.t. } \mathcal{F}, OUTF(a) \neq \emptyset, INF(b) \neq \emptyset\}$ 
9     $FP := \{p \cup q \mid p \cap q \neq \emptyset, p \in R, q \in \Delta R, \mathcal{WS} \models p \cup q\};$ 
10   forall  $(a, b) \in FA$  do
11     forall  $p_1 \in OUTF(a), p_2 \in INF(b)$  s.t.  $(a, b) \notin p_1 \cup p_2$  and  $(p_1, p_2) \notin connected\_by(a, b)$  do begin
12        $q := p_1 \cup p_2 \cup \{(a, b)\};$ 
13       if  $\mathcal{WS} \models q$  then begin
14          $FP := FP \cup \{q\};$ 
15          $INBORDER(q) := ComputeInBorder(b, p_1, p_2);$ 
16          $OUTBORDER(q) := ComputeOutBorder(a, p_1, p_2);$ 
17         forall  $a \in INBORDER(q)$  do  $INFP(a) := INFP(a) \cup \{q\};$ 
18         forall  $a \in OUTBORDER(q)$  do  $OUTFP(a) := OUTFP(a) \cup \{q\};$ 
19          $connected\_by(a, b) := connected\_by(a, b) \cup \{(p_1, p_2)\};$ 
20       end
21     end
22    $\Delta R := \{p \in FP \mid p \text{ is frequent w.r.t. } \mathcal{F}\};$ 
23    $R := R \cup \Delta R;$ 
24 until  $\Delta R = \emptyset;$ 
25 return  $R;$ 

```

Function $ComputeInBorder(b, p_1, p_2);$
if $|InDegree_{p_1 \cup p_2}(b)| < InDegree(b) - 1$ **then** $INBORDER := \{b\}$ **else** $INBORDER := \emptyset;$
forall $c \in (INBORDER(p_1) \cup INBORDER(p_2)) - \{b\}$ **do**
 if $|InDegree_{p_1 \cup p_2}(c)| < InDegree(b)$ **then** $INBORDER := INBORDER \cup \{c\};$
return $INBORDER;$

Function $ComputeOutBorder(a, p_1, p_2);$
if $|OutDegree_{p_1 \cup p_2}(a)| < OutDegree(a) - 1$ **then** $OUTBORDER := \{a\}$ **else** $OUTBORDER := \emptyset;$
forall $c \in (OUTBORDER(p_1) \cup OUTBORDER(p_2)) - \{a\}$ **do**
 if $(|OutDegree_{p_1 \cup p_2}(c)| < OutDegree(a))$ **then** $OUTBORDER := OUTBORDER \cup \{c\};$
return $OUTBORDER;$

Fig. 7. Algorithm $c\text{-find}(\mathcal{F}, \mathcal{WS})$.

Proof. It is easy to see that $c\text{-find}$ considers all the candidate patterns considered by $w\text{-find}$. This is a trivial consequence of Lemma 3 and of the observation that $c\text{-find}$ degenerates in $w\text{-find}$ each time it considers elementary patterns in R . This also entails $N_c \leq N_w$. However, in general, $\mathcal{C} = \mathcal{W}$ does not hold. Indeed, let us consider the situation in which there are four patterns p_1, p_2, p_3 , and p_4 and the patterns $p_1 \cup p_2$ and $p_2 \cup p_4$ are frequent, but the pattern $p_1 \cup p_3$ is not. Assume also that patterns $p_1 \cup p_2$ and $p_3 \cup p_4$ are connected by means of an edge (a, b) . In such a case, $c\text{-find}$ would generate the (unfrequent) candidate pattern $p_1 \cup p_2 \cup p_3 \cup p_4 \cup \{(a, b)\}$, but $w\text{-find}$ would not. \square

As a consequence of the above statements, the two algorithms can be considered as viable alternatives, and a preference can be carried out only by considering the particular structure of the workflow schema that had generated the instances. The next section will also provide a discussion and a comparison between the two algorithms.

5 EXPERIMENTS AND DISCUSSION

In this section, we study the behavior of the algorithms $w\text{-find}$ and $c\text{-find}$ by evaluating both their performance and their scalability. As shown in the previous section, the algorithms are sound and complete with regard to the set of frequent w -patterns. Nevertheless, in principle, the number

of candidate w -patterns generated could be prohibitively high, thus making the algorithms unfeasible on complex workflow schemas. Moreover, we also compare the performance of our implementations with regard to several existing techniques for computing frequent itemsets adapted to the particular applicative domain.

In our experiments, we mainly use synthetic data. Synthetic data generation can be tuned according to:

1. the size of \mathcal{WS} ,
2. the size of \mathcal{F} ,
3. the size $|L|$ of the frequent weak patterns in \mathcal{F} , and
4. the probability p_{\subseteq} of choosing a E^{\subseteq} -arc.

The ideas adopted in building the generator for synthetic data are essentially inspired by [21].

5.1 Performance of $w\text{-find}$

In a first set of experiments, we tested the $w\text{-find}$ algorithm by first considering some fixed workflow schemas and generating synthesized workflow instances. In particular, the nondeterministic choices in the executions are performed according to a binomial distribution with mean p_{\subseteq} . Frequent instances are forced into the system by replicating some instances (in which some variations were randomly performed) according to $|L|$. Fig. 8a reports the number of operations (matching of a pattern with an instance) for

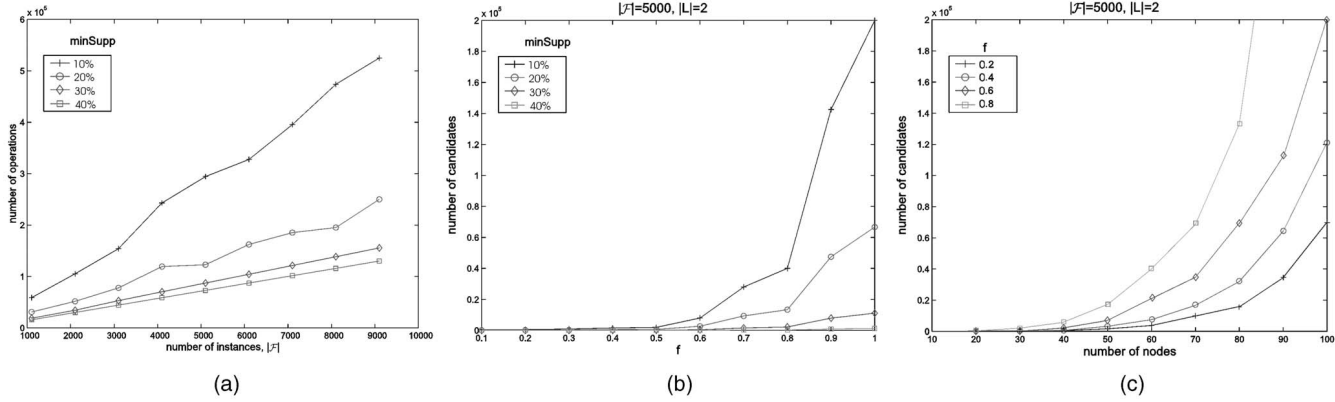


Fig. 8. w -find performance. (a) Number of operations with regard to $|\mathcal{F}|$. (b) Number of candidates with regard to f for different \minSupp values. (c) Number of candidates with regard to number of nodes for randomly generated workflow schemas.

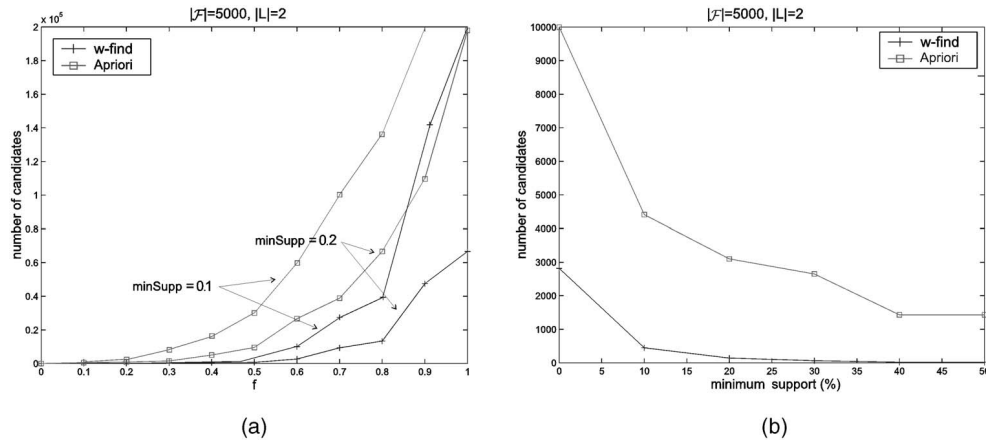


Fig. 9. w -find versus Apriori. (a) Number of candidates with regard to f . (b) Number of candidates for different \minSupp values.

increasing values of $|\mathcal{F}|$. The figure shows that the algorithm scales linearly in the size of the input (for different supports).

In a second set of experiments, we randomly generate the workflow schemas to test the efficiency of the approach with regard to the structure of the workflow. To this aim, we fix $|\mathcal{F}|$ and generate workflow instances according to the randomly generated schema. The actual number of nodes and arcs is chosen by picking from a Poisson distribution with fixed mean value. In order to evaluate the contribution of the complexity of workflow schemas, we exploit the factor $f = \frac{|E^\subseteq|}{|E|}$, which represents the degree of potential nondeterminism within a workflow schema. Intuitively, workflow schemas exhibiting $f \simeq 0$ produce instances with a small number of candidate w -patterns. Conversely, workflow schemas exhibiting $f \simeq 1$ produce instances with a large number of candidate w -patterns. Fig. 8b shows the behavior of w -find when f ranges between 0 (no nondeterminism) and 1 (full nondeterminism) for different values of \minSupp values. It is interesting to observe that even for significantly higher values of f (real workflow schema are expected to have a degree of nondeterminism less than 0.5), the smart way of searching the search space reduces drastically the number of candidates being generated. Finally, in Fig. 8c, we report the number of candidates at the varying of the number of nodes for different values of f . It is worth noting the exponential behavior due to the combinatorial explosion of the search space.

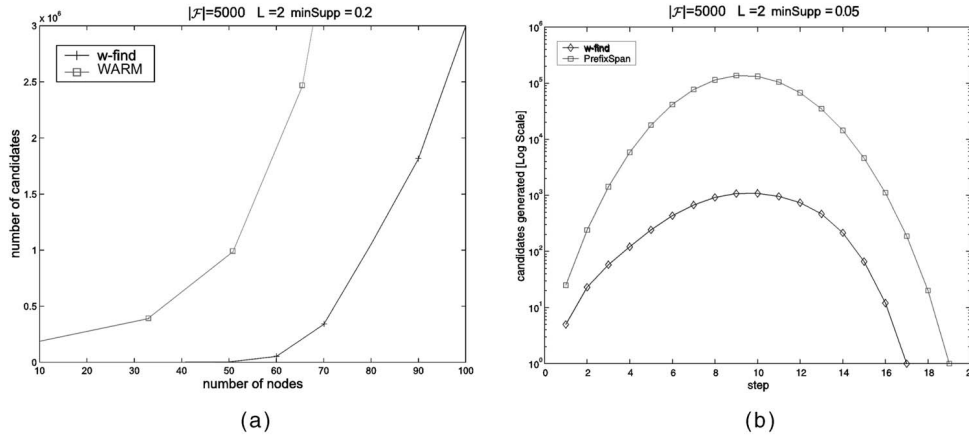
5.2 Comparing w -find and Apriori

We consider an implementation of the *Apriori* algorithm, which only computes frequent itemsets of edges in E^\subseteq . Such an approach is significant for analyzing the performance overhead suffered by traditional frequent-pattern mining methods, which typically can be easily adapted to mine workflow instances, but are not tuned to take into account domain information about the workflow schema. We perform several experiments comparing the performance of the *Apriori* approach with the ones of w -find on increasing values of $|\mathcal{F}|$ and \minSupp . For a data set of instances generated as said before with regard to the workflow schema of Fig. 1, the comparison is reported in Fig. 9b.

As expected, w -find outperforms *Apriori* by an order of magnitude. This is mainly due to the fact that, contrary to w -find, in the *Apriori* implementation, arcs in E^\subseteq are combined without taking into account the information provided by the workflow schema.

Fig. 9a shows instead the behavior of both *Apriori* and w -find when f ranges between 0 (no nondeterminism) and 1 (full nondeterminism). Again, *Apriori* is outperformed by w -find. Notice that for small values of f , both the algorithms produce a small number of candidates; however, in this situation, w -find still performs significantly better than *Apriori* for small \minSupp values (e.g., 0.1). In fact, for lower values of \minSupp , the number of candidates increments significantly and, hence, the focused strategy of w -find leads to a significant gain. However, we point out that the adaptation of *Apriori* tested here might be a viable solution

warmodekey(instance(-I)).	rmode(startarc(+I, #N)).	type(startnode(pic, obj)).
talking(3).	rmode(endnode(+I, #N).	type(endnode(pic, obj)).
usepacks(0).	rmode(andarc(+I, #S, #D)).	type(andarc(pic, obj, obj)).
minfreq(0.2).	rmode(xorarc(+I, #S, #D)).	type(xorarc(pic, obj, obj)).
typed ₁ language(yes).	rmode(optarc(+I, #S, #D)).	type(optarc(pic, obj, obj)).
	rmode(arc(+I, #S, #D)).	type(arc(pic, obj, obj)).
	rmode(node(+I, #N)).	type(node(pic, obj)).

Fig. 10. The setting file used in the *WARMR* approach.Fig. 11. (a) Comparison of *w-find* with *WARM* over a fixed workflow schema. (b) Comparison of *w-find* with *PrefixSpan*.

for the mining of “nearly deterministic” workflows, if we are, moreover, interested in very frequent patterns ($\text{minSupp} > 0.2$).

5.3 Comparing *w-find* with *WARMR* and *PrefixSpan*

A possible further approach to consider is the *WARMR* algorithm, devised in [26], which allows an explicit formalization of domain knowledge (like, for example, the connectivity information provided by the workflow schema) which can be directly exploited by the algorithm. The setting file that we have used is reported in Fig. 10 (predicates to be mined are on the right).

The results of the comparison are shown in Fig. 11a, where we report the correlation between the number of candidate patterns and the number of the nodes in the workflow schema at the varying of f .

In the evaluation of the algorithm, we also have made some comparison with regard to methods for mining sequential pattern. However, a workflow is not a sequence; nonetheless, we can assume to represent each instance as a sequential pattern by considering the ordering of execution of each activity. For example, the instance reported in Fig. 3 can be described by the sequence $s_1 = \langle S, c, (efg), (e_2l), (4, 5), j_2, j, l, (mno), A \rangle$ if we assume that each activity requires the same amount of time to be executed. Note that we grouped all the activities that we assume to be executed at the same time. Conversely, if the activity g requires more time, a possible ordering of executions associated to the same instance is $s_2 = \langle S, c, e, (e_2), (4, 5), j_2, j, g, l, (mno), A \rangle$. Thus, s_1 and s_2 are distinct sequences associated to the same instance. It follows that any sequential pattern algorithm can be used for extracting frequent instances, but it cannot be complete in the sense that some frequent instances will not be mined since the sequences associated to the executions are possibly quite different (and, hence, infrequent). In our testing, we compared *w-find* with the *PrefixSpan* algorithm [4], but in order to achieve a finer analysis, we assume each activity to require the same time to be executed; essentially,

PrefixSpan has been applied on the sequences constructed from each instance by performing a breadth-first search, starting from the initial activity.

The results are reported in Fig. 11b, where we report the number of candidates generated at each stage of the computation for a fixed workflow schema. Again, this experiment is significative only for understanding the advantage of a more focused method and is not a comparison on “pure” sequences where *PrefixSpan* is expected to outperform both *w-find* and *c-find*. In fact, we can see that the more elaborate and domain dependent way of searching patterns in the lattice leads to a smaller number of patterns to be generated.

5.4 Comparing *w-find* and *c-find*

Finally, we report the experimental results of the comparisons between *w-find* and *c-find*. In a first set of experiments, we fixed a value of f (0.7) and generated 5,000 random instances. In Fig. 12a, we report the number of candidates generated over such instances at the varying of the minimum support. It is interesting to observe that *w-find* performs better than *c-find*, especially for lower values of minSupp .

For a second set of experiments, we fixed $\text{minSupp} = 0.2$, and we made the comparison at the varying of f . This second set of experiments, whose results are reported in Fig. 12b confirmed the quality of *w-find* of generating fewer candidates, for every type of workflow (regardless of the degree of nondeterminism).

The factor that may instead lead to a preference of *c-find* is in the number of steps performed. Let us consider Fig. 13, which reports the number of candidates generated at the different steps of the algorithm (the scale is logarithmic). Here, the behavior of *c-find* is somehow dual to that of *w-find* (as reported in Fig. 11b). Indeed, *c-find* at each successive step generates more candidates than in the previous one and this leads the process to converge quickly. Conversely, *w-find* after a certain number of steps dramatically reduces

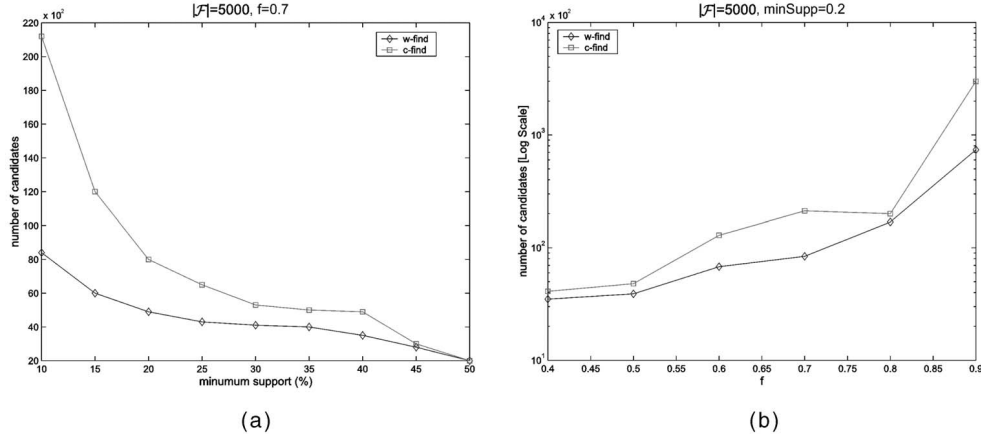


Fig. 12. Comparison of *w-find* with *c-find*. (a) Number of candidates for different *minSupp* values. (b) Number of candidates at the varying of the nondeterminism.

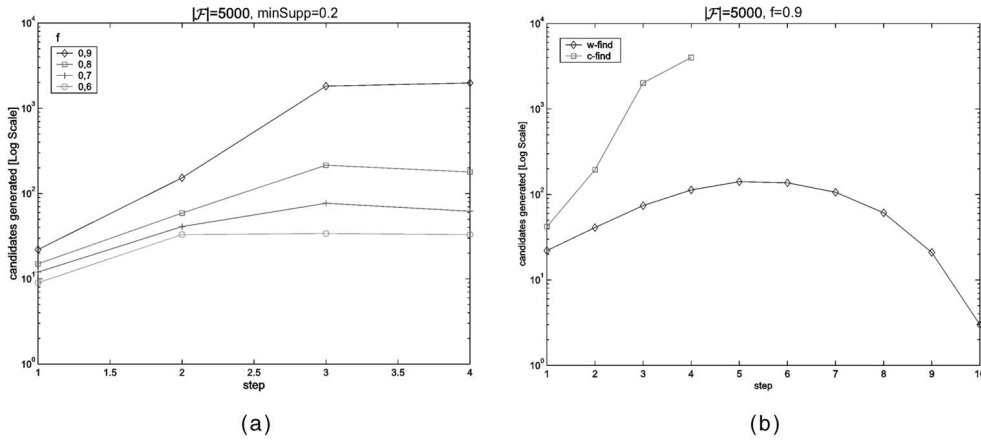


Fig. 13. (a) Number of candidates generated at the different steps. (b) Comparison of *w-find* with *c-find*.

the number of new frequent patterns discovered and, hence, it requires more iterations. A more direct comparison is reported in Fig. 13b, from which it is evident that the faster rate of convergence of *c-find* is paid with a bigger number of candidates generated. Since the number of steps coincides with the number of scans of the database, in the case of huge databases, *c-find* may be convenient.

More generally, *c-find* is expected to exhibit better performance than *w-find* with *dense* workflow databases. More specifically, a set of workflow instances is dense if the number of expected frequent patterns is large with regard to the size of the workflow. If a database of instances exhibits this peculiarity, the number of candidate patterns to be generated is likely to be of the same order of magnitude as the number of frequent patterns (that is, the number of unfrequent patterns is small with regard to the set of frequent ones). In such a case, both *c-find* and *w-find* are expected to compute (almost) the same set of candidates. However, the look-ahead strategy of the *c-find* algorithm guarantees a faster convergence rate. Clearly, more efficient extensions could be devised to the proposed algorithms for dense databases in order to avoid candidate generation (e.g., in the style of [3]).

6 CONCLUSIONS

We have introduced the problem of mining frequent instances of workflow schemas, motivated by the aim of

providing facilities for the system administrator to monitor the actual behavior of the workflow system in order to predict the “most probable” workflow executions. We have shown that the use of mining techniques is justified by the fact that even “simple” reachability problems are intractable.

We have proposed two novel graph mining algorithms specialized to deal with constraints imposed by the structures of workflow schemes and instances and we have studied their properties, both theoretically and experimentally, by showing that they represent an effective means of investigating some inherent properties of the executions of a given schema.

Following our approach, future research might develop more elaborated algorithms that are able to deal with more expressive modeling features which have been not considered in our formal framework. For instance, a valuable on-going extension is dealing with supporting cyclic instances by integrating our techniques with well-known approaches for mining periodic patterns (see, e.g., [31]).

ACKNOWLEDGMENTS

This work was partially supported by Project SP1 on “Reti Internet: efficienza, integrazione e sicurezza” and Project SFIDA-PMI, both funded by MIUR.

REFERENCES

- [1] R. Agrawal and R. Srikant, "Mining Sequential Patterns," *Proc. 11th Int'l Conf. Data Eng.*, pp. 3-14, 1995.
- [2] M. Zaki, "Efficiently Mining Frequent Trees in a Forest," *Proc. Eighth Int'l Conf. Knowledge Discovery and Data Mining*, pp. 71-80, 2002.
- [3] J. Han, J. Pei, and Y. Yi, "Mining Frequent Patterns without Candidate Generation," *Proc. Int'l ACM Conf. Management of Data*, pp. 1-12, 2000.
- [4] J. Pei, J. Han, B. Mortazavi-Asl, H. Pinto, Q. Chen, U. Dayal, and M. Hsu, "Prefixspan: Mining Sequential Patterns by Prefix-Projected Growth," *Proc. IEEE Int'l Conf. Data Eng.*, pp. 215-224, 2001.
- [5] A. Inokuchi, T. Washi, and H. Motoda, "An Apriori-Based Algorithm for Mining Frequent Substructures from Graph Data," *Proc. Fourth European Conf. Principles of Data Mining and Knowledge Discovery*, pp. 13-23, 2000.
- [6] X. Yan and J. Han, "gSpan: Graph-Based Substructure Pattern Mining," *Proc. IEEE Int'l Conf. Data Mining*, an extended version appeared as UIUC-CS Technical Report R-2002-2296, 2001.
- [7] M. Kuramochi and G. Karypis, "Frequent Subgraph Discovery," *Proc. IEEE Int'l Conf. Data Mining*, pp. 313-320, 2001.
- [8] X. Yan and J. Han, "CloseGraph: Mining Closed Frequent Graph Patterns," *Proc. ACM Int'l Conf. Knowledge Discovery and Data Mining*, pp. 286-295, 2003.
- [9] P. Senkul, M. Kifer, and I.H. Toroslu, "A Logical Framework for Scheduling Workflows under Resource Allocation Constraints," *Proc. 28th Int'l Conf. Very Large Data Bases*, pp. 694-702, 2002.
- [10] H. Schuldt, G. Alonso, C. Beeri, and H. Schek, "Atomicity and Isolation for Transactional Processes," *ACM Trans. Database Systems*, vol. 27, no. 1, pp. 63-116, 2002.
- [11] A. Bonner, "Workflow, Transactions, and Datalog," *Proc. 18th ACM Symp. Principles of Database Systems*, pp. 294-305, 1999.
- [12] H. Davulcu, M. Kifer, C.R. Ramakrishnan, and I.V. Ramakrishnan, "Logic Based Modeling and Analysis of Workflows," *Proc. 17th ACM Symp. Principles of Database Systems*, pp. 25-33, 1998.
- [13] D. Wodtke and G. Weikum, "A Formal Foundation for Distributed Workflow Execution Based on State Charts," *Proc. Sixth Int'l Conf. Database Theory*, pp. 230-246, 1997.
- [14] D. Wodtkeand, J. Weissenfels, G. Weikum, and A. Dittrich, "The Mentor Project: Steps Towards Enterprise-Wide Workflow Management," *Proc. IEEE Int'l Conf. Data Eng.*, pp. 556-565, 1996.
- [15] G. Kappel, P. Lang, S. Rausch-Schott, and W. Retschitzagger, "Workflow Management Based on Object, Rules, and Roles," *IEEE Data Eng. Bull.*, vol. 18, no. 1, pp. 11-18, 1995.
- [16] M.P. Sing, "Semantical Considerations on Workflows: An Algebra for Intertask Dependencies," *Proc. Int'l Workshop Database Programming Languages*, pp. 6-8, 1995.
- [17] W.M.P. van der Aalst, "The Application of Petri Nets to Workflow Management," *J. Circuits, Systems, and Computers*, vol. 8, no. 1, pp. 21-66, 1998.
- [18] W.M.P. van der Aalst, A. Hirsenschall, and H.M.W. Verbeek, "An Alternative Way to Analyze Workflow Graphs," *Proc. 14th Int'l Conf. Advanced Information Systems Eng.*, pp. 534-552, 2002.
- [19] R. Agrawal, D. Gunopulos, and F. Leymann, "Mining Process Models from Workflow Logs," *Proc. Sixth Int'l Conf. Extending Database Technology*, pp. 469-483, 1998.
- [20] W.M.P. van der Aalst, B.F. van Dongen, J. Herbst, L. Maruster, G. Schimm, and A.J.M.M. Weijters, "Workflow Mining: A Survey of Issues and Approaches," *Data and Knowledge Eng.*, vol. 47, no. 3, pp. 237-267, 2003.
- [21] R. Agrawal and R. Srikant, "Fast Algorithms for Mining Association Rules," *Proc. 20th Int'l Conf. Very Large Databases*, 1994.
- [22] J. Pei, J. Han, H. Lu, S. Nishio, S. Tang, and D. Yang, "H-Mine: Hyper-Structure Mining Of Frequent Patterns in Large Databases," *Proc. IEEE Int'l Conf. Data Mining*, pp. 441-448, 2001.
- [23] P. Koksai, S.N. Arpinar, and A. Dogac, "Workflow History Management," *SIGMOD Record Archive*, vol. 27, no. 1, pp. 67-75, 1998.
- [24] W.M.P. van der Aalst and K.M. van Hee, *Workflow Management: Models, Methods, and Systems*. MIT Press, 2002.
- [25] J.E. Cook and A.L. Wolf, "Automating Process Discovery through Event-Data Analysis," *Proc. 17th Int'l Conf. Software Eng.*, pp. 73-82, 1995.
- [26] L. Dehaspe and H. Toivonen, "Discovery of Frequent DATALOG Patterns," *Data Mining and Knowledge Discovery*, vol. 3, no. 1, pp. 7-36, 1999.
- [27] D. Georgakopoulos, M. Hornick, and A. Sheth, "An Overview of Workflow Management: From Process Modeling to Workflow Automation Infrastructure," *Distributed and Parallel Databases*, vol. 3, no. 2, pp. 119-153, 1995.
- [28] M.R. Garey and D.S. Johnson, *Computers and Intractability: A Guide to the Theory of NP-completeness*. New York: Freeman and Comp., 1979.
- [29] N.D. Jones and T. Laaser, "Complete Problems for Deterministic Polynomial Time," *Theoretical Computer Science*, vol. 3, pp. 105-117, 1977.
- [30] W.M.P. van der Aalst, A.H.M. ter Hofstede, B. Kiepuszewski, and A.P. Barros, "Advanced Workflow Patterns," *Proc. Seventh Int'l Conf. Cooperative Information Systems*, pp. 18-29, 2000.
- [31] J. Yang, W. Wang, and P.S. Yu, "Mining Asynchronous Periodic Patterns in Time Series Data," *Proc. Sixth ACM SIGKDD Int'l Conf. Knowledge Discovery and Data Mining*, pp. 275-279, 2000.



Gianluigi Greco received the Laurea degree in computer science engineering from the University of Calabria, Italy, in 2000. Currently, he is an assistant professor of computer science in the Department of Mathematics at the University of Calabria. His main research interests are in the areas of databases and artificial intelligence with a focus on database theory, knowledge representation, nonmonotonic reasoning, computational complexity, and deductive databases.



Antonella Guzzo received the Laurea degree in engineering from the University of Calabria. She has been a PhD student in computer engineering at the DEIS Department of the University of Calabria since 2001. Her research interests include workflow management systems, knowledge representation, and data mining.



Giuseppe Manco graduated summa cum laude in computer science in 1994 and received the PhD degree in computer science from the University of Pisa. He is currently a senior researcher at the Institute of High Performance Computing and Networks (ICAR-CNR) of the National Research Council of Italy and a contract professor at University of Calabria, Italy. He has been contract researcher at the CNUCE Institute in Pisa, Italy, and a visiting fellow at the CWI Institute in Amsterdam, Netherlands. His current research interests include deductive databases, knowledge discovery and data mining, Web databases, and semistructured data.



Domenico Saccà received the Doctoral degree in engineering from the University of Rome in 1975. Since 1987, he has been a full professor of computer engineering at the University of Calabria and, since 1995, he has also been the director of the CNR (the Italian National Research Council) Research Institute ICAR (Institute for High Performance Computing and Networking). In the past, he was a visiting scientist at the IBM Laboratory of San Jose, at the Computer Science Department of UCLA, and at the ICSI Institute of Berkeley; furthermore, he was a scientific consultant of MCC, Austin, and manager of the Research Division of CRAI. His current research interests focus on advanced issues on database such as scheme integration in data warehousing, compressed representation of data-cubes, workflow and process mining, and logic-based database query languages. He has been a member of the program committees of several international conferences, director of international schools and seminars, and leader of many national and international research projects. He is a senior member of the IEEE Computer Society.