

A Context-aware Miner for Medical Processes

L. Canensi (1), G. Leonardi (2), S. Montani (2), P. Terenziani (2)

(1) Department of Computer Science, Università di Torino, Italy

(2) DISIT, Computer Science Institute, Università del Piemonte Orientale,
Alessandria, Italy

Abstract. Medical process mining is gaining much attention in recent years, but the available mining algorithms can hardly cope with medical application peculiarities, that require to properly *contextualize* process patterns. Indeed, most approaches loose the connection between a mined pattern and the relevant portion of the input event log, and can have a limited precision, i.e., they can mine incorrect paths, never appearing in the input log traces. These issues can be very harmful in medical applications, where it is vital that mining results are reliable as much as possible, and properly reference the contextual information, in order to facilitate the work of physicians and hospital managers in guaranteeing the highest quality of service to patients.

In this paper we propose a novel approach to medical process mining that operates in a **context-aware** fashion. We show on a set of critical examples how our algorithm is able to cope with the issues sketched above. The process model we mine can then be adopted to support efficient and flexible **trace retrieval**, as described in [1]. Indeed, the model can be used as an indexing structure, well suited to quickly retrieve traces corresponding to the pattern being looked for. In the future, we plan to test the approach on a real-world medical dataset, taken from the stroke patient management domain.

1 Introduction

Process mining describes a family of a-posteriori analysis techniques that exploit the so-called *event log*, which records information about the sequences (*traces* henceforth) of activities executed at a given organization. The most relevant and widely used process mining technique is discovery; process discovery takes as an input the event log and produces a process model, without using any a-priori information.

Despite the complexity of healthcare, medical process mining is a research field which is gaining attention in recent years (see, e.g., [2–4]), also as a means for assessing the correct application of clinical guideline directives in practice. This application domain indeed presents particular challenges and issues [5], mostly related to the fact that different types of patients exhibit different characteristics, that the patient’s state dynamically evolves (and is influenced by the medical activities executed on her/him), and that different hospital settings may have

different resource constraints. All these peculiarities lead to the need to properly *contextualize* medical processes and/or process patterns.

The currently available process mining solutions, ranging from commercial tools (offered, e.g., by Fujitsu Ltd and Celonis), to the open-source framework ProM [6] (developed at the Eindhoven University of Technology), which represents the state of the art in process mining research, are not tailored to medical applications, and do not take into account contextualization needs.

Specifically, despite some differences, many current algorithms show important similarities (see Section 2), and have common limitations: (1) they learn “context-free” patterns of processes; (2) they can mine paths that do not correspond to any input trace in the log (i.e., they can have a limited *precision* [7]); (3) they do not explicitly relate the mined patterns to the log (in the sense that there is no explicit correspondence between the mined patterns, and the traces in the log “supporting” them).

Such limitations are quite relevant in general, and very relevant in the medical domain. Concerning limitation (1), it is well known that, e.g., the same (set of) activities may produce different effects on patients, depending on the context (e.g., on the activities previously performed on the patients themselves). The impact of limitation (2) is obvious and dramatic: the mined model is the input for quality assessment procedures, such as verification of conformance with respect to clinical guidelines, or performance measurement and bottleneck detection. All these procedures will provide an unreliable output, if played on an unreliable input (in the sense that the model may exhibit a path that never appears in any input trace). However, surprisingly, limited precision is a common limitation of many current miners (see Section 2.) Limitation (3) is less critical, but still significant. Indeed, maintaining an explicit link between the mined patterns and the input traces matching such patterns, can be important not only to characterize contexts, but also to provide physicians with an evidence of the learned output.

In this paper, we propose an innovative approach that supports **context-aware** process mining, and overcomes all the above limitations. Our mining technique provides a process model in the form of a tree, called the “log-tree”. The log-tree perfectly adheres to the input traces (since all its paths correspond to at least one trace in the log), and explicitly relates the mined patterns to the log. Interestingly, in this way it can be exploited as an index, to efficiently retrieve all the traces that contain the pattern of interest, along the lines of the work in [1], where we describe a **trace retrieval** approach for business process management operational support¹.

¹ This paper is however focused on the mining technique, and on the description of the log-tree structure. The interested reader can find additional details about trace retrieval in [1].

2 A critical analysis of current approaches

As discussed in the Introduction, several different miners have been developed in the literature. However, many existing miners, including *alpha miner* [8], *fuzzy miner* [9], *heuristic miner* [10] and the very recent *inductive tree miner* [11], show interesting commonalities.

In the following, for the sake of clarity and brevity, we will mainly refer to one miner to illustrate the common characteristics of these literature approaches. Specifically, we will concentrate on heuristic miner. Indeed, heuristic miner can abstract from exceptional behavior and noise and, therefore, is suitable for many real-life logs, including medical ones.

Most mining algorithms are based on a common assumption, which we call **uniqueness assumption**: each event is unique in the output, so that each event appears at most once in the output of the process miner².

As a second commonality, the mining methodology is typically based on two steps (corresponding to steps 1 and 3 in Algorithm 1 below, which abstractly characterizes the behavior of heuristic miner on a set T of traces): (i) a **de-structuring step**, in which immediate precedence between pairs of events is detected in the input traces, and (ii) a **re-structuring step**, in which patterns of events are reconstructed, by combining the immediate precedence relations mined in the de-structuring step.

ALGORITHM 1: HM pseudocode

```

1 function HM(T) ;
2 (1) for each pair of events  $A$  and  $B$  do
3   | search T to detect immediate precedence  $A \rightarrow B$  and  $B \rightarrow A$ 
4 end
5 (2) Look for cycles of length one (same event repeated) or two (pair of events
   repeated);
6 (3) for each triplet of events  $A, B$ , and  $C$ , such that  $A \rightarrow B$  and  $A \rightarrow C$  do
7   | use formulae [10] to discriminate whether  $B$  and  $C$  are in AND (both of them
   | must be executed after  $A$ , in any order) or in XOR (exactly one of them must be
   | executed after  $A$ ), and construct the output
8 end
9 (4) Look for “long-distance dependencies” between events, and add them to the
   output.
```

For the sake of brevity, here we cannot comment in detail the algorithm, and the formulae. We just focus on the main critical issues.

C1 Given the uniqueness assumption, the de-structuring step (step 1 in Algorithm 1) evaluates the immediate precedence relations between events A

² Actually, ProM’s *region-based miner* [12] is able to relax this assumption in some cases; however, it does not maintain or exploit contextual information; moreover, its precision can be very low, making this algorithm less suitable for our purposes.

- and B looking at sequences “AB” and “BA” in the input traces, regardless of their position in the traces, and, thus, regardless of the context.
- C2 Given the fact that the de-structuring step ignores the context, and that the re-structuring step (step 3 in Algorithm 1) does not take into account the traces in the log, also the re-structuring step does not consider the context at all.
- C3 The uniqueness assumption imposes severe constraints on the re-structuring step.

The following examples illustrate the critical impact of issues C1, C2 and C3. The examples have been processed using heuristic miner, and provided as Petri Nets, which are commonly assumed as an incontestable formalism to represent (the semantics of) processes.

To simplify the presentation, with no loss of generality, we suppose that each trace is prefixed with a distinguished starting symbol (*) and postfixed with another distinguished ending symbol (#).

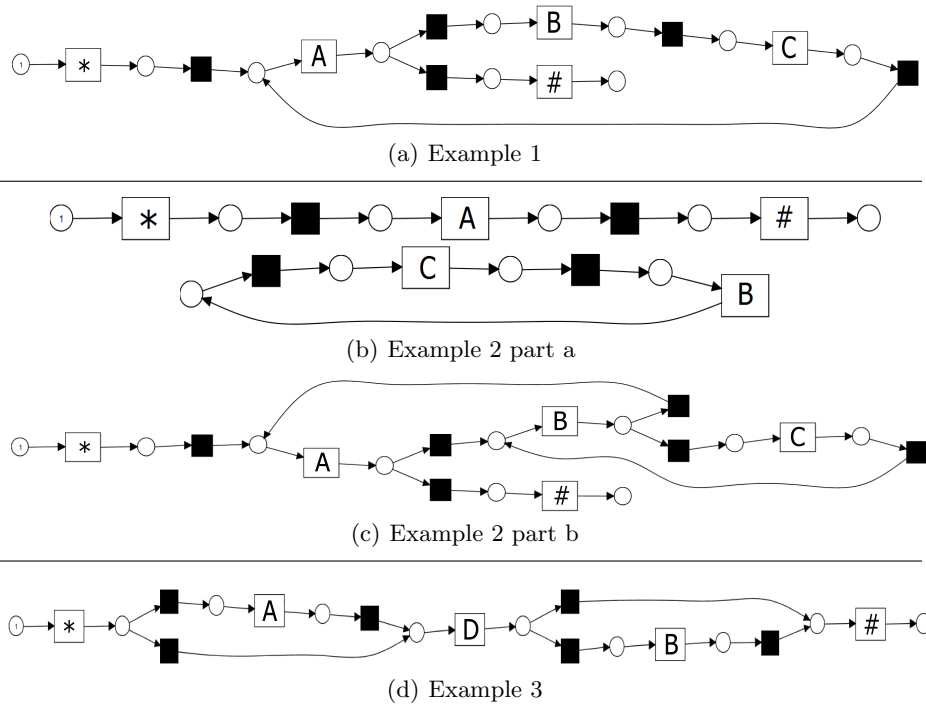


Fig. 1. Models mined by heuristic miner referring to a set of critical examples

Ex.1 Log content: 1000 equal traces: *ABCA#
 Although all the traces are equal, existing literature approaches cannot learn

the (unique!) pattern. In this example, this is mainly due to the fact that the uniqueness assumption causes problems in the re-structuring phase (critical issue C3): the two occurrences of the event A in the traces must correspond to a unique transition in the output Petri Net. Thus, a cycle (returning from C to A) must be introduced in the output. Notably, the output Petri Net also admits patterns like $*A\#$ or $*ABCABCA\#$, which are not present in any of the input traces (see Figure 1(a)).

The pattern shown in this example represents one of the typical care processes for patients suffering from acute diseases. For instance, consider the diagnostic process of patients suffering from cerebral ischemia. After the symptom onset, the patient arrives at the emergency ward of the hospital. According to the latest medical guidelines for the treatment of stroke [13,14], it is recommended that the patient undergoes as soon as possible a computed tomography (CT), for: (1) the differential diagnosis between ischemic and hemorrhagic stroke and other cerebrovascular diseases, and (2) the identification of possible early signs of ischemic brain suffering. This action is indicated in Ex. 1 as action A = CT Execution. Subsequently, the patient should be evaluated by a trained neurologist, generating action B = Neurological Evaluation. Among the various additional investigations, which may be required depending on the type of patient, the guidelines suggest to always perform an electrocardiogram (ECG). In Ex. 1, this is referred to as action C = ECG Execution. Before transferring the patient to the Stroke Unit for further treatment, the guidelines recommend the repetition of the CT (within 48 hours), in order to obtain a better diagnostic and prognostic evaluation. The action A = CT Execution is therefore performed both as the first action of this process, and as the last one, thus generating the pattern of actions of this example.

In a similar manner, each of the following examples, presented only in an abstract way due to lack of space, introduces situations which can be concretely instantiated as real clinical processes, or parts of them.

Ex.2 Log content: 1000 equal traces: $*ABCBA\#$

This example highlights the limitations of “context-free” approaches. Consider, in particular, critical issue C1 above: since the precedence relations are searched without considering the context (and on the basis of the uniqueness assumption), there is no way to decide whether A precedes B (100% of traces, considering the first part of the traces) or B precedes A (100% of traces, but considering the last part of the traces), and analogously for the relation between B and C. As a consequence, two separate Petri Nets are learned: the first contains only the pattern $*A\#$, while the second contains a loop composed by C and B (see Figure 1(b)). By using heuristic miner without the “all event connected option”, a different model is learned (see Figure 1(c)), where it is possible to replay the input trace, but many other never observed behaviors can be generated as well.

Ex.3 Log content: 500 traces $*AD\#$ and 500 traces $*DB\#$

Once again, the combination of a “context-free” analysis with the uniqueness as-

sumption causes undesired effects. In this example, two different patterns should be mined. However, during the de-structuring phase, no distinction is made between D in the first type of traces (i.e., D in the context of being preceded by *A) and D in the second type of traces (i.e., D in the context of being preceded by *). Thus, the re-structuring phase (see critical issue C2) generates a “merging” between the two different patterns (due to the uniqueness of D), providing the Petri Net in Figure 1(d) as output. Notably, besides the correct patterns *AD# and *DB#, the output Petri Net also models the patterns *D# and ADB#, which do not correspond to any trace in the input log.

3 Context-aware process mining

In order to overcome the limitations illustrated in the previous section, we propose an innovative approach to (medical) process mining, which is based on quite a different philosophy: in our methodology, process mining heavily takes into account contextual information, both in the data structure, and in the mining algorithm. Our approach is based on three main ideas:

- in the data structure, we **remove the uniqueness assumption**: the same event may appear more than once in the output model, to represent the fact that it may occur in different contexts;
- in the data structure, and in the mining algorithm, we explicitly maintain the context, i.e., the set of log traces that *support* a given path in the mined model;
- in the mining algorithm, we take advantage of the ordering of events in the log traces (which represents, indeed, the context in which each single event occurs) to directly mine the output model, **without distinguishing between a de-structuring and a re-structuring phase**.

In the following, we present our approach and its application to the previously discussed critical examples.

3.1 Data structure and algorithm

Our mining algorithm takes in input a log (set of traces). The log is represented by a matrix, in which each row is a trace. It outputs the mined process as a tree, called a “log-tree”, in which nodes represent events (i.e., activities), and arcs represent a precedence relation between them. More precisely, in our model, each node is represented as a pair $\langle P, T \rangle$:

- P denotes a (possibly unary) set of events; events in the same node are in *AND* relation, or, more properly, may occur in any order (with respect to each other). Note that, in such a way, each path from the root node of the tree to a given node N denotes a set of possible process patterns (called *support patterns* of N henceforth), obtained by following the order represented by

- the arcs in the path to visit the tree, and ordering in each possible way the events in each node (for instance, the path $\{A, B\} \rightarrow \{C\}$ represents the support patterns “ABC” and “BAC”);
- T represents the context, i.e., a set of references to all and only those traces in the log which exactly match one of the patterns in P (called *support traces* henceforth).

Our mining algorithm (see Algorithm 2 below) builds the log-tree described as above.

ALGORITHM 2: Mining pseudocode

```

1 function Build-Tree(index, < P, T >);
2 nextP ← getNext(index+1, T);
3 if nextP not empty then
4   nextEvents ← XORvsAND (nextP, T);
5   for each node < P', T' > ∈ nextEvents do
6     AppendSon(< P', T' >, < P, T >);
7     Build-Tree(index + |P'|, < P', T' >);
8   end
9 end

```

Build-Tree in Algorithm 2 takes in input a variable *index*, representing a given position in the traces (i.e., a column in the input matrix), and a node. Initially, it is called on the first position, and on the root of the tree (which is a “dummy” node, corresponding to the * event; thus, initially, index=0, P=* and T is the set of all the traces). The function *getNext* simply inspects the traces in T to find all possible next events (in the context T). On the basis of the current context (*support traces*) T , the function *XORvsAND* applies the formulae described in appendix A to identify which events are in AND and which are in XOR relation. The output of such a function is a set of nodes $\langle P', T' \rangle$, one for each maximal set of events to be AND-ed. Note that, for each one of such sets P' , the corresponding set T' of *support traces* is also computed, on the basis of the current context T . Finally, each new node is appended in the output tree (function *AppendSon*), and *Build-Tree* is recursively applied to each node (with the parameter *index* properly set).

A-posteriori, we create a dummy node #, and connect all the leaves to it; in this way, the log-tree can also be seen as a basic process model.

Notably, our mining algorithm explicitly manages the context, focusing at each step on the proper *support traces*, and always taking into account the global ordering of events in the traces to build the graph. As a consequence, differently from the heuristic miner algorithm shown in Algorithm 1, we do not need any additional step to cope with long-distance dependencies (we have them “for free”). Analogously, we “naturally” cope with cycles by simply unfolding them. Thus, our algorithm already directly copes with cycles of *any length* (and no

additional ad-hoc procedure for cycles is needed, differently from the heuristic miner algorithm).

3.2 Examples

In Figure 2, we show the output of our miner, applied to the examples Ex.1-Ex.3 discussed in Section 2. Support traces are not reported in the figure, but are part of the output itself. As it can be observed, the critical situations discussed in Section 2 are all correctly managed by our approach. It is also worth noting that our process graphs could be easily converted into Petri Nets.

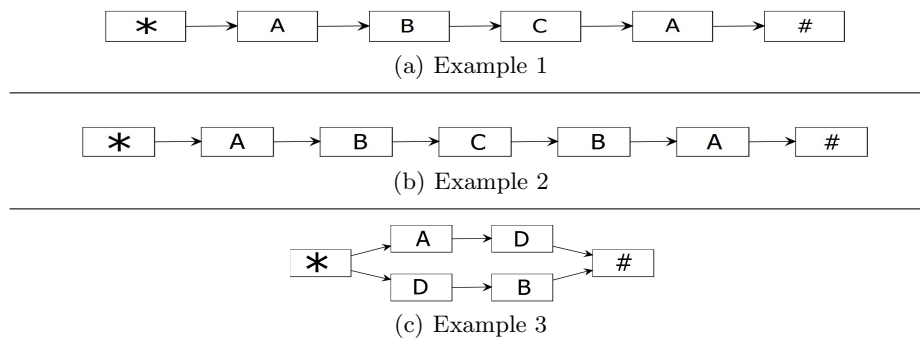


Fig. 2. Models mined by our miner referring to the examples of Section 2

Quite naturally, the log-tree is not only a model of the process at hand, but also an index of the event log content, since, by means of the *support traces*, it explicitly relates the mined patterns to the log. Indeed, as described in [1], we are exploiting the log-tree structure to speed up trace retrieval, in a tool for operational support we are currently completing and testing.

4 Conclusions

In this paper, we have introduced a novel process mining algorithm, able to overcome the limitations of many current approaches available in the literature. Specifically, our algorithm:

- learns “context-aware” patterns of processes;
- has a high precision, since it provides patterns that always correspond to input traces in the log;
- explicitly relates the mined patterns to the traces in the log “supporting” them.

Our approach properly deals with critical situations that may occur in practical application domains, as illustrated by the examples in Section 2. These characteristics make the algorithm particularly well-suited for medical applications, where it is vital that mining results are reliable as much as possible, in order to facilitate the work of physicians and hospital managers in guaranteeing the highest quality of service to patients.

The mined process model can also be exploited to allow efficient trace retrieval, as described in [1]. Indeed, the model we mine maintains an explicit link between mined patterns and the input traces supporting them, and can thus be used as an indexing structure, well suited to quickly retrieve traces corresponding to the pattern at hand.

In the future, we plan to extensively test the overall approach on a real-world medical dataset, taken from the stroke patient management domain.

References

1. A. Bottrighi, L. Canensi, G. Leonardi, S. Montani, and P. Terenziani. Trace retrieval for business process operational support. *Expert Syst. Appl.*, 55:212–221, 2016.
2. E. Rojas, J. Munoz-Gama, M. Sepulveda, and D. Capurro. Process mining in healthcare: A literature review. *Journal of Biomedical Informatics*, 61:224 – 236, 2016.
3. R. Mans, H. Schonenberg, G. Leonardi, S. Panzarasa, A. Cavallini, S. Quaglini, and et al. Process mining techniques: an application to stroke care. In S. Andersen, G. O. Klein, S. Schulz, and J. Aarts, editors, *Proc. MIE*, volume 136 of *Studies in Health Technology and Informatics*, pages 573–578. IOS Press, Amsterdam, 2008.
4. L. Perimal-Lewis, S. Qin, C. Thompson, and P. Hakendorf. Gaining insight from patient journey data using a process-oriented analysis approach. In K. Butler-Henderson and K. Gray, editors, *Proc. Workshop on Health Informatics and Knowledge Management (HIKM)*, volume 129 of *Conferences in Research and Practice in Information Technology*, page 5966. Australian Computer Society, 2012.
5. R. Mans, W. van der Aalst, R. Vanwersch, and A. Moleman. Process mining in healthcare: Data challenges when answering frequently posed questions. In R. Lenz, S. Miksch, M. Peleg, M. Reichert, D. Riaño, and A. ten Teije, editors, *ProHealth/KR4HC*, volume 7738 of *Lecture Notes in Computer Science*, pages 140–153. Springer, Berlin, 2013.
6. B. van Dongen, A. Alves De Medeiros, H. Verbeek, A. Weijters, and W. Van der Aalst. The proM framework: a new era in process mining tool support. In G. Ciardo and P. Darondeau, editors, *Knowledge Management and its Integrative Elements*, pages 444–454. Springer, Berlin, 2005.
7. J. Buijs, B. van Dongen, and W. van der Aalst. On the role of fitness, precision, generalization and simplicity in process discovery. In *On the Move to Meaningful Internet Systems: OTM 2012*, pages 305–322. Springer, 2012.
8. W. Van der Aalst and B. van Dongen. Discovering workflow performance models from timed logs. In Y. Han, S. Tai, and D. Wikarski, editors, *Proc. International Conference on Engineering and Deployment of Cooperative Information Systems (EDCIS 2002)*, volume 2480 of *Lecture Notes in Computer Science*, pages 45–63. Springer, Berlin, 2002.

9. C. Günther and W. van der Aalst. Fuzzy mining - adaptive process simplification based on multi-perspective metrics. In G. Alonso, P. Dadam, and M. Rosemann, editors, *Business Process Management, 5th International Conference, BPM 2007, Brisbane, Australia, September 24-28, 2007, Proceedings*, volume 4714 of *Lecture Notes in Computer Science*, pages 328–343. Springer, 2007.
10. A. Weijters, W. Van der Aalst, and A. Alves de Medeiros. *Process Mining with the Heuristic Miner Algorithm, WP 166*. Eindhoven University of Technology, Eindhoven, 2006.
11. S. Leemans, D. Fahland, and W. van der Aalst. Discovering block-structured process models from event logs containing infrequent behaviour. In N. Lohmann, M. Song, and P. Wohed, editors, *Business Process Management Workshops - BPM International Workshops, Beijing, China, August 26, 2013, Revised Papers*, volume 171 of *Lecture Notes in Business Information Processing*, pages 66–78. Springer, 2013.
12. J. Carmona, J. Cortadella, and M. Kishinevsky. A Region-Based Algorithm for Discovering Petri Nets from Event Logs. In *Business Process Management (BPM2008)*, pages 358–373, 2008.
13. G. Carlucci D. Inzitari. Italian stroke guidelines (spread): evidence and clinical practice. *Neurological Sciences*, 27:S225–S227, 1996.
14. <http://www.iso-spread.it>. The ISO-SPREAD clinical guideline group web page.

A Xor vs And Formulae

In Algorithm 2, after finding the set of successors $nextP$ of the considered set of events P , we focus on the discovery of the relation between them. In order to do this, we calculate the *dependency frequency* between every event pairs $\langle A, B \rangle$ in $nextP \times nextP$:

$$A \rightarrow B = \frac{1}{2} \left(\frac{|A > B|}{\sum_{X \in Act_T} |A > X|} + \frac{|A > B|}{\sum_{Y \in Act_T} |Y > B|} \right) \quad (1)$$

where $|A > B|$ is the number of traces in which A is immediately followed by B, and $|A > X|$ is the number of traces in which A is immediately followed by some event X, $|Y > B|$ is the number of traces in which B is immediately preceded by some event Y. After evaluating the dependency frequency value $A \rightarrow B$ and $B \rightarrow A$, we can have the following possible situations:

- if both the values are below a given threshold, this means that A and B rarely appear in the same trace, therefore they are in XOR relation;
- if $A \rightarrow B$ is above the threshold and $B \rightarrow A$ is below, then A precedes B, and, viceversa, if $B \rightarrow A$ is above the threshold and $A \rightarrow B$ is below, then B precedes A;
- if both the values are above the threshold, then A and B are in the any order relation.