

# Improving Process Discovery Results by Filtering Outliers Using Conditional Behavioural Probabilities

Mohammadreza Fani Sani<sup>(✉)</sup>, Sebastiaan J. van Zelst,  
and Wil M. P. van der Aalst

Department of Mathematics and Computer Science,  
Eindhoven University of Technology,  
P.O. Box 513, 5600 MB Eindhoven, The Netherlands  
{m.fani.sani,s.j.v.zelst,w.m.p.v.d.aalst}@tue.nl

**Abstract.** Process discovery, one of the key challenges in process mining, aims at discovering process models from process execution data stored in event logs. Most discovery algorithms assume that all data in an event log conform to correct execution of the process, and hence, incorporate all behaviour in their resulting process model. However, in real event logs, noise and irrelevant infrequent behaviour are often present. Incorporating such behaviour results in complex, incomprehensible process models concealing the correct and/or relevant behaviour of the underlying process. In this paper, we propose a novel general purpose filtering method that exploits observed conditional probabilities between sequences of activities. The method has been implemented in both the ProM toolkit and the RapidProM framework. We evaluate our approach using real and synthetic event data. The results show that the proposed method accurately removes irrelevant behaviour and, indeed, improves process discovery results.

**Keywords:** Process mining · Process discovery · Noise filtering  
Outlier detection

## 1 Introduction

*Process mining* is a research discipline that positioned at the intersection of data driven methods like machine learning and data mining and Business Process Modeling (BPM) [1]. There are three types of process mining; *process discovery*, *conformance checking* and *process enhancement*. Process discovery aims at discovering process models from event logs. Conformance checking aims at assessing to what degree a process model and event log conform to one another in terms of behaviour. Finally, process enhancement aims at improving process model quality by enriching them with information gained from the event log.

Within process mining/process identification projects, process discovery is often used to quickly get insights regarding the process under study [1]. A business process analyst simply applies a process discovery algorithm on the

extracted event log and analyzes its result. Most process discovery algorithms assume that event logs represent accurate behaviour. So, they are designed to incorporate all of the event log's behaviour in their resulting process model as much as possible.

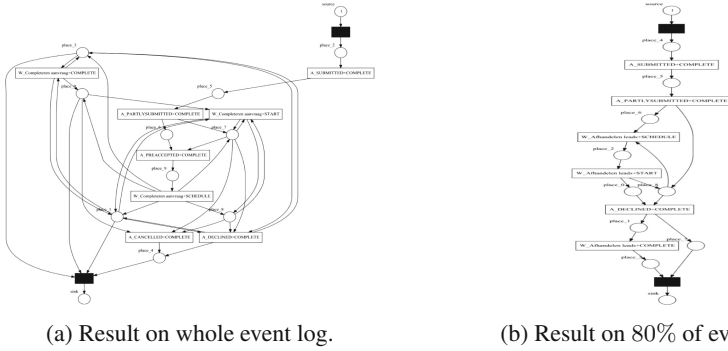
Real event logs contain both *noise* and *infrequent behaviour* [2]. In general, noise refers to behaviour that does not conform to the process specification and/or its correct execution. Examples of noise are, amongst others, incomplete logging of process behaviour, duplicated logging of events and faulty execution of the process. Infrequent behaviour relates to behaviour that is may be supposed to happen, yet, in very exceptional cases of the process. For example, additional checks may be required when a loan request exceeds \$10.000.000. Incorporating either noise and/or infrequent behaviour results in complex, incomprehensible process models concealing the correct and/or relevant behaviour of the underlying process. As such, when using process discovery for the purpose of process identification, we are often unable to gain any actionable knowledge by applying process discovery algorithms directly.

In this paper, we focus on improving process discovery results by applying general purpose event log filtering, i.e. filtering the event log prior to applying any arbitrary process discovery algorithm. Distinguishing between noise and infrequent behaviour is a challenging task and is outside the scope of this paper. Hence, we consider both noise and infrequent behaviour as *outliers* and aim at identifying and removing such outliers from event logs. We propose a generic filtering approach based on conditional probabilities between *sequences of activities*. The approach identifies whether certain activities are likely to happen based on a number of its preceding activities. Using the ProM (<http://promtools.org>) [3] based extension of RapidMiner (<http://rapidminer.com>), i.e. RapidProM [4], we study the effectiveness of our approach, using synthetic and real event data. The results of our experiments show that our approach adequately identifies and removes outliers, and, as a consequence increases the overall quality of process discovery results. Additionally we show that our filtering method outperforms other general purpose filtering techniques in the process mining domain.

The remainder of this paper is structured as follows. Section 2 motivates the need for general purpose event log filtering methods. In Section 3, we discuss related work and after that, in Sect. 4, we explain our proposed method. Details of the evaluation and corresponding results are given in Sect. 5. Finally, Sect. 6 concludes the paper and presents some future work in this domain.

## 2 Motivation

An interpretable process model helps business process analysts to understand what is going on in an event data. However, often process discovery algorithms return results that are complicated and not understandable, because of outliers within the event logs used. In Fig. 1 we show how the application of filtering greatly reduces the complexity in a *real event log*, i.e. the event log of the *Business Process Intelligence Challenge 2012 (BPIC 2012)*. Figure 1a shows a process



**Fig. 1.** Process models discovered by applying the ILP Miner [5] on the BPIC 2012 log.

model discovered using the ILP Miner of [5]<sup>1</sup> for this event log, whereas Fig. 1b shows the result of applying the same process discovery algorithm on 80% of the most frequent original behaviour.

In process mining, two quality measures are defined for measuring the behavioural quality of process models, i.e. fitness and precision [6]. Fitness computes how much behaviour in the event log is also described by the process model. On the other hand, precision measures the amount of behaviour described by the model that is also present in the event log. The fitness values of Figs. 1a and 1b are 0.57 and 0.46 whereas their precision values are 0.68 and 1.0 respectively. This means that the model in Fig. 1a describes more behaviour that is also presented in the event log, however, in order to do this it is greatly under-fitting. Hence, it allows for much more behaviour compared to the model in Fig. 1b. As a consequence, the model in Fig. 1a is more complex and ambiguous. However, arbitrarily removing behaviour based on frequency is too ad-hoc and does not work when there is a lot of variety present within an event log. Therefore, we need more advanced filtering methods that take into account and exploit the actual behaviour described by the event log. Clearly, by incorporating all possible behaviour, the model in Fig. 1a is overly complex and conceals the dominant/main-stream behaviour of the underlying process. The process model in Fig. 1b, on the other hand, is much simpler while still covering 80% of the observed behaviour allowed by the other model. At the same time, the model describes less behaviour compared to the model in Fig. 1a, and is therefore less under-fitting. Thus, by removing 20% of behaviour, we obtain a simpler model that still accurately describes the underlying process.

### 3 Related Work

In recent years, many process discovery algorithms have been proposed [7–12]. The first algorithms were designed to incorporate all behaviour in the event

<sup>1</sup> HybridILPMiner package in ProM.

log [7, 11, 12]. However, later more algorithms have more recently been extended to be able to handle outliers [13, 14]. However, these extended filtering techniques are tailored towards the internal working of the corresponding algorithm and hence do not work as a general purpose filtering technique. Other algorithms [8, 10], are designed to cope with noisy and infrequent behaviour, however, these algorithms do not result in process models with a clear semantics. Most of commercial process mining tools using these algorithms and their filtering are based on just frequency of activities and their direct relations. Moreover, the filters are relatively ad-hoc and require significant user input.

In this paper, we propose to *separate concerns*, and thus develop a novel, general purpose filtering technique that pre-processes event logs. In such way, any process discovery algorithm is able to benefit from effective identification and removal of outlier behaviour. In the remainder of this section, we focus on techniques developed for general purpose filtering in the process mining domain.

Separating outliers from event logs and focusing just on them rather than all behaviour also has been studied [15], however, a detailed treatment of such work is outside the scope of this paper.

The vast majority of process mining research has an accompanying implementation in the process mining toolkit ProM. Most work on general purpose event log filtering concerns ad-hoc filtering implementations within ProM. Many of these implementations are useful when we aim at using specific subsets of traces/events of an event log instead of the whole event log. In Table 1, the main filtering plugins are listed, accompanied by a brief description of their applications and methods. All plugins take an event log as an input and return a *filtered*

**Table 1.** Overview of filtering plugins in ProM

Plug-in	Applications	Main method
Filter log using simple heuristics	Helpful for removing traces and activities based on frequency of events or the presence of certain start/end events	Frequency/position of events
Filter log on event/trace attributes	Useful when we want to just keep events/traces with specific attribute values	Attribute values
Dotted chart	Allows us to visually select specific traces in event logs (usually base on a time frame)	Time window
Transition systems miner	Helpful to project traces/events on specific transitions and/or states	Frequency of transitions
Filter log using prefix-closed language	Allows us to remove events from traces	Rule based

*event log* as an output. Moreover, they need some form of domain knowledge to work properly. In addition, typically the user needs to set one or more (complex) settings. However, they do not support generic outlier detection, i.e. in cases where we possess no or little domain knowledge.

Surprisingly, little research has been done in the field of general purpose filtering. In [16] a graph-based outlier detection method is proposed to detect inaccurate data in an event log. In [17] a method is proposed that detects non-fitting behaviour based on a given reference model and then repair event log. As we want to improve process discovery results and, in general, we do not have a reference model, this method is not useful for general purpose filtering. In [18] the authors propose to provide training traces to the PRISM algorithm [19] which returns rules for detecting outliers. However, in real event logs, providing a set of training traces that covers all possible outliers is impractical. Also, in this method, deciding about level of filtering is impossible.

The most relevant research in the area of general purpose log filtering is the work in [20]. The authors propose to construct an Anomaly Free Automaton (AFA) based on the whole event log and a given threshold. Subsequently, all events that do not fit the AFA are removed from the filtered event log. Filtering event logs using AFA indeed allows us to detect and remove noisy and/or infrequent behaviour. However, the technique does not allow us to detect all types of outliers like incomplete traces, i.e. traces that fit the AFA perfectly yet do not terminate properly. Incorporation of such behaviour can still lead to infeasible process discovery results.

## 4 Filtering with Conditional Behavioural Probabilities

As indicated in Sect. 3, the most filtering approaches are not suitable for process discovery because they need additional information like reference model or a set of outlier traces. Furthermore, the AFA filter, which is the most suitable general purpose event log filter, has trouble identifying irrelevant infrequent behaviour. Therefore, we present a general purpose filtering method that is able to deal with all types of outliers. The main purpose of the filter is to identify the likelihood of the occurrence of an activity, based on its surrounding behaviour, e.g. how likely is it that activity  $a$  follows the sequence of activities  $\langle b, c \rangle$ . To detect such likelihood it uses the conditional probability of activity occurrences, given a sequence of activities. As we just consider a sample of behaviour in the underlying process, i.e. event log, all computed probabilities are estimation of behaviour that are really happened. Prior to presenting the filtering method, we present some basic notations used throughout the paper.

### 4.1 Basic Notation and Definitions

Given a set  $X$ , a multiset  $M$  over  $X$  is a function  $M: X \rightarrow \mathbb{N}_{\geq 0}$ , i.e. it allows certain elements of  $X$  to appear multiple times. We write a multiset as  $M = [e_1^{k_1}, e_2^{k_2}, \dots, e_n^{k_n}]$ , where for  $1 \leq i \leq n$  we have  $M(e_i) = k_i$  with  $k_i \in \mathbb{N}$ . If  $k_i = 1$ ,

we omit its superscript, and if for some  $e \in X$  we have  $M(e) = 0$ , we omit it from the multiset notation. Also,  $M = []$  is an empty multi set if  $\forall e \in X, M(e) = 0$ . We let  $\bar{M} = \{e \in X \mid M(e) > 0\}$ , i.e.  $\bar{M} \subseteq X$ . The set of all possible multisets over a set  $X$  is written as  $\mathcal{M}$ .

Let  $\mathcal{A}$  denote the set of all possible activities and let  $\mathcal{A}^*$  denote the set of all finite sequences over  $\mathcal{A}$ . A finite sequence  $\sigma$  of length  $n$  over  $\mathcal{A}$  is a function  $\sigma: \{1, 2, \dots, n\} \rightarrow \mathcal{A}$ , alternatively written as  $\sigma = \langle a_1, a_2, \dots, a_n \rangle$  where  $a_i = \sigma(i)$  for  $1 \leq i \leq n$ . The empty sequence is written as  $\epsilon$ . Concatenation of sequences  $\sigma$  and  $\sigma'$  is written as  $\sigma \cdot \sigma'$ . We let  $hd: \mathcal{A}^* \times \mathbb{N}_{\geq 0} \rightarrow \mathcal{A}^*$  with, given some  $\sigma \in \mathcal{A}^*$  and  $k \leq |\sigma|$ ,  $hd(\sigma, k) = \langle a_1, a_2, \dots, a_k \rangle$ , i.e., the sequence of the first  $k$  elements of  $\sigma$ . Note that  $hd(\sigma, 0) = \epsilon$ . Symmetrically  $tl: \mathcal{A}^* \times \mathbb{N}_{\geq 0} \rightarrow \mathcal{A}^*$  is defined as  $tl(\sigma, k) = \langle a_{n-k+1}, a_{n-k+2}, \dots, a_n \rangle$ , i.e., the sequence of the last  $k$  elements of  $\sigma$ . Again,  $tl(\sigma, 0) = \epsilon$ . Finally, sequence  $\sigma' = \langle a'_1, a'_2, \dots, a'_k \rangle$  is a subsequence of sequence  $\sigma$  if and only if we are able to write  $\sigma$  as  $\sigma_1 \cdot \langle a'_1, a'_2, \dots, a'_k \rangle \cdot \sigma_2$ , where both  $\sigma_1$  and  $\sigma_2$  are allowed to be  $\epsilon$ , i.e.  $\sigma$  is a subsequence of itself.

Event logs describe sequences of executed business process activities, typically in context of some case, e.g. a customer or some order-id. The execution of an activity in context of a case is referred to as an *event*. A sequence of events is referred to as a *trace*. Thus, it is possible that multiple traces describe the same sequence of activities, yet, each trace itself contains different events. An example event log, adopted from [1], is presented in Table 2.

**Table 2.** Fragment of a fictional event log (each line corresponds to an event).

Case-id	Activity	Resource	Time-stamp
...	...	...	...
1	register request (a)	Sara	2017-04-08:08.10
1	examine thoroughly (b)	Ali	2017-04-08:09.17
2	register request (a)	Sara	2017-04-08:10.14
2	check ticket (d)	William	2017-04-08:10.23
1	check ticket (d)	William	2017-04-08:10.53
2	examine causally (b)	Ava	2017-04-08:11.13
1	reject request (h)	Ava	2017-04-08:13.05
...	...	...	...

Consider all activities related to *Case-id 1*. Sara registers a request, after which Ali examines it thoroughly. William checks the ticket after which Ava examine causally and reject the request. The execution of an activity in context of a business process is referred to as an *event*. A sequence of events, e.g. the sequence of events related to case 1, is referred to as a *trace*.

**Definition 1 (Trace, Variant, Event Log).** Let  $\mathcal{A}$  be a set of activities. An event log is a multiset of sequences over  $\mathcal{A}$ , i.e.  $L \in \mathcal{M}(\mathcal{A}^*)$ .  $\sigma \in \mathcal{A}^*$  is a trace in  $L$  and  $\sigma \in \bar{L}$  is a variant.

Observe that each  $\sigma \in \bar{L}$  describes a *trace-variant* whereas  $L(\sigma)$  describes how many traces of the form  $\sigma$  are present.

**Definition 2 (Subsequence Frequency).** Let  $L$  be an event log over a set of activities  $\mathcal{A}$  and let  $\sigma' \in \mathcal{A}^*$ . The subsequence frequency of  $\sigma'$  w.r.t  $L$ , written as  $freq(\sigma', L)$ , denotes the number of times  $\sigma'$  occurs as a subsequence of any trace present in  $L$ .

Given a simple example event log  $L = [\langle a, b, c, d \rangle^5, \langle a, c, b, d \rangle^3]$ , we have  $freq(\langle a \rangle, L) = freq(\epsilon, L) = 8$ ,  $freq(\langle a, b \rangle, L) = 5$ , etc.

**Definition 3 (Conditional Occurrence Probability).** Let  $L$  be an event log over a set of activities  $\mathcal{A}$  and let  $\sigma'$  be a subsequence. Given some  $a \in \mathcal{A}$ , the conditional probability of occurrence of activity  $a$ , given  $\sigma'$  and  $L$ , i.e.  $COP(a, \sigma', L)$  is defined as:

$$COP(a, \sigma', L) = \begin{cases} \frac{freq(\sigma' \cdot \langle a \rangle, L)}{freq(\sigma', L)} & \text{if } freq(\sigma', L) \neq 0 \\ 0 & \text{otherwise} \end{cases}$$

Clearly, the value of any  $COP(a, \sigma, L)$  is a real number in  $[0, 1]$ . A high value of  $COP(a, \sigma, L)$  implies that after the occurrence of  $\sigma'$ , it is very probable that activity  $a$  occurs. For example,  $COP(a, \sigma, L) = 1$  implies that if  $\sigma$  occurs,  $a$  always happens directly after it. Based on the previously used simple event log, we have  $COP(b, \langle a \rangle, L) = \frac{5}{8}$ .

## 4.2 Outlier Detection

We aim to exploit conditional probabilities present within event logs for the purpose of filtering event logs. Conceptually, after a given subsequence, activities that have a particularly low  $COP$ -value are unlikely to have happened and therefore their occurrence may be seen as outlier. However, to account for dependencies between activities and previously occurred activities at larger distances, we compute  $COP$ -values for subsequences of increasing length.

In our proposed method, for each  $i \in \{1, 2, \dots, k\}$  we construct a  $COP$ -Matrix. Assume there are a total of  $m$  unique subsequences with length  $1 \leq l \leq k$  in an event log. A  $COP$ -Matrix  $\mathbf{A}_{COP}^l$  for length  $l$  is simply an  $m \times |\mathcal{A}|$ -matrix, where  $\mathbf{A}_{COP}^l(\sigma', a) = COP(a, \sigma', L)$ .

We additionally compute conditional probabilities for *start* and *end* subsequences relatively. We let  $\mathbf{A}_S^l$  denote a matrix describing the occurrence probability matrix of all subsequences  $\sigma' = hd(\sigma)$  with  $|\sigma'| = l$  for  $\sigma \in L$ . We are able to compute such probability by dividing the number of traces that start with  $\sigma'$  over the total number of traces in the log. Similarly we define  $\mathbf{A}_E^l$  denote a matrix describing the conditional probability matrix of all subsequences  $\sigma' = tl(\sigma)$  with  $|\sigma'| = l$  for  $\sigma \in L$  that is equal to  $a = \epsilon$  in  $\mathbf{A}_{COP}^l$ . By doing so, we be able to handle outliers which occur in the start and the end parts of trace.

Given our different  $COP$ -Matrix, and a user-defined threshold  $\kappa$ , we identify each entry  $\mathbf{A}^l(\sigma', a) < \kappa$  as an outlier. The pseudo-code of detecting outliers is

present in Algorithm 1. In this fashion, it is possible to detect outliers that occur in start, middle or end part of traces. There are two ways to handle detected outliers. We are able to simply remove the corresponding event from the trace, i.e. *event-level filtering*, or, remove the trace as a whole, i.e. *trace-level filtering*. However, removing an improbable event in a trace may make the trace to have more outlier behaviour. Hence, we just focus on *trace-level filtering*.

---

**Algorithm 1.** Outlier Detection Algorithm
 

---

```

procedure OUTLIERDETECTION( $L, k, \kappa$ )
  Computing Probabilities:
  for ( $l \leq k$ ) do
    Build  $\mathbf{A}_{COP}^l, \mathbf{A}_E^l$  and  $\mathbf{A}_S^l$ 
  FilteredEventLog  $\leftarrow$  EmptyEventLog
  Filtering:
  for (each Trace  $\sigma$  in the  $L$ ) do
    Outlier  $\leftarrow$  false
    for ( $l = 1 : k$ ) do
      for (each subsequence  $\sigma'$  with length  $l$  and its following activity  $a$ ) do
        Find corresponding  $COP(a, \sigma', L)$  in  $\mathbf{A}_{COP}^l, \mathbf{A}_E^l$  and  $\mathbf{A}_S^l$ 
        if ( $\kappa > COP(a, \sigma', L)$ ) then
          Outlier  $\leftarrow$  true
    if (Outlier = false) then
      FilteredEventLog  $\leftarrow$  Add(FilteredExceptionLog,  $\sigma$ )
  return FilteredExceptionLog
  
```

---

With increasing value of  $k$  (maximum length of subsequences), the complexity of the filtering method increases. The number of different strings we can generate over  $\mathcal{A}$  with length  $k$  is  $(|\mathcal{A}|)^k$  and total possible subsequences for some  $k$ :  $\sum_{i=1}^k \mathcal{A}^i$  where  $|\mathcal{A}|$  is the number of activities in the  $L$ . However, there is no need to compute  $COPs$  of all possible subsequences. For subsequences with length  $k + 1$ , it is sufficient to just consider  $\sigma'.\langle a \rangle$  in level  $k$  that  $\kappa \leq COP(a, \sigma', L)$ . For example, if at  $k = 1$   $COP(c, \langle b \rangle) \leq \kappa$ , there is no need to consider  $\langle b, c \rangle$  as a subsequent at  $k = 2$ . Even the  $COP(a, \langle b, c \rangle)$  be higher than  $\kappa$ .

### 4.3 Implementation

To be able to combine the proposed filtering method with any process discovery algorithm, we implemented the *Matrix Filter* plugin (*MF*) in the **ProM** framework<sup>2</sup>. The plugin takes an event log as input and outputs a filtered event log. Furthermore, the user is able to specify threshold  $\kappa$  and whether event-level or trace-level filtering needs to be applied. The maximum subsequence length to be considered also needs to be specified.

In addition, to apply our proposed method on various event logs with different filtering thresholds and applying different process discovery algorithms with different parameters, we ported the *Matrix Filter* (*MF*) plugin to **RapidProM**. **RapidProM** is an extension of **RapidMiner** that combines scientific workflows [21] with a range of (**ProM**-based) process mining algorithms.

---

<sup>2</sup> MatrixFilter plugin [svn.win.tue.nl/repos/prom/Packages/LogFiltering](http://svn.win.tue.nl/repos/prom/Packages/LogFiltering).



## 5 Evaluation

To evaluate the usefulness of filtering outliers using our method, we have conducted several experiments using both synthetic and real event data. The purpose of these experiments is to answer the following questions:

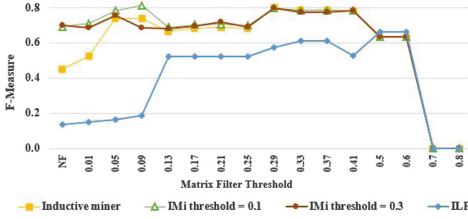
1. Does *MF* help process discovery algorithms to return more precise models?
2. How does the performance of *MF* compare to *AFA* filtering method?

To evaluate discovered process models, we use fitness and precision (introduced in Page 2). There is a trade off between these measures [22]. Sometimes, putting aside little behaviour cause a few decrease in fitness value, however more increasing in the precision value. To make a balance between fitness and precision, we use the F-Measures metric that combines fitness and precision:  $\frac{2 \times \text{Precision} \times \text{Fitness}}{\text{Precision} + \text{Fitness}}$ . Also, filtering time and process model discovery time in milliseconds have been measured. Note that in all experiments, filtered event logs are only used in the process discovery part. Computing the F-Measure for all process models is done using the corresponding raw, unfiltered event logs. Furthermore, we only consider subsequences with length  $k$  in  $[0, 2]$ .

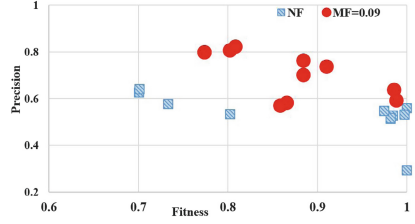
In the first experiment we investigate the effect of changing the  $\kappa$  in the *MF* threshold on the F-Measure w.r.t. different process discovery algorithms. We use the Inductive Miner [12] (IM) and the ILP Miner (ILP) [11]. Additionally we assess the interaction between our filtering technique and integrated filtering within the Inductive Miner, i.e. we use the IMi variant [13] with noise thresholds 0.1 and 0.3. We apply these algorithms and filtering methods on the BPIC2012 log. The results for this experiment are shown in Fig. 2.

In this figure, each line corresponds to a discovery algorithm. The  $x$ -axis represents the threshold level of *MF*, the  $y$ -axis represents the corresponding F-Measure. Hence, for each technique, the data point  $x = 0$  corresponds to not applying behavioural conditional probability filtering. We thus aim at finding out whether there exist regions of threshold values for which we are able to increase the F-measure when applying filtering. The F-measure of *IM* on this event log without using *MF* is 0.45. However, using the proposed filter increases the F-Measure of the discovered model to 0.80. Even for *IMi*, which uses an embedded filter, the *MF* increases the F-Measure from 0.69 and 0.7 to 0.81. As the ILP miner is more sensitive to outliers, *MF* helps more and its enhancement for this algorithm is higher. However, with increasing the threshold of *MF*, all the traces in the filtered event log are removed and the fitness and F-Measure of the discovered model will equal to 0. The best result, i.e. an F-measure of 0.81, is achieved by IMi with threshold 0.1 and *MF* threshold of 0.09.

To illustrate the effect of filtering on the discovered process models, in Fig. 3, we apply *IMi* with 11 internal thresholds ranging from 0.0 to 0.5 on the raw BPIC2012 and the filtered event log using *MF* with threshold 0.09. Here, each circle or square correspond to fitness and precision values related to one discovered model. A circle is related to applying *MF*, whereas squares relate to using the raw event log. As the results show, *MF* usually causes a little decrease in fitness value, yet yields an increase in precision value. The average of F-Measures



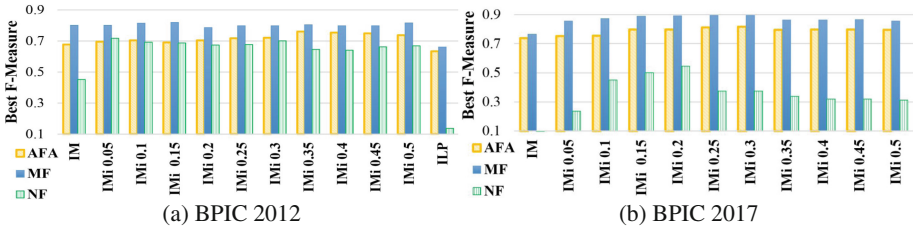
**Fig. 2.** Applying process discovery algorithms on the BPIC2012 log with different *MF* thresholds.



**Fig. 3.** Comparing process models discovered by 11 noise thresholds on the BPIC 2012 log with/without filtering.

when applying no filtering is 0.66 versus 0.77 in case of *MF* (with threshold 0.09). Thus, Figs. 2 and 3 indicate that *MF* improve process discovery results, i.e. the process models have an overall higher F-Measure.

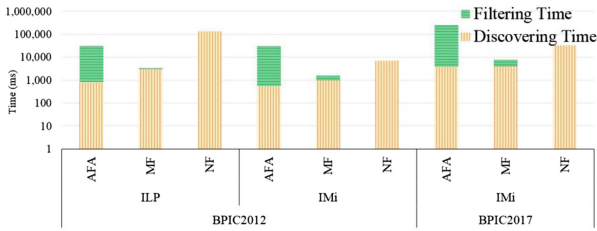
In next experiment, using the BPIC2012 and BPIC2017 event log we additionally assess what the maximal obtainable level of F-measure is for different process discovery algorithms, using different levels of internal filtering. We computed F-measures based on the unfiltered event log, and, maximized the F-measure result for both *MF* and *AFA*. With a workflow in the *RapidMiner*, for both filtering methods we filtered the event log using 40 different thresholds. The results are presented in Fig. 4. This figure shows *MF* allows us to discover process models with higher F-Measures.



**Fig. 4.** Effect of filtering on best F-Measure of discovered models.

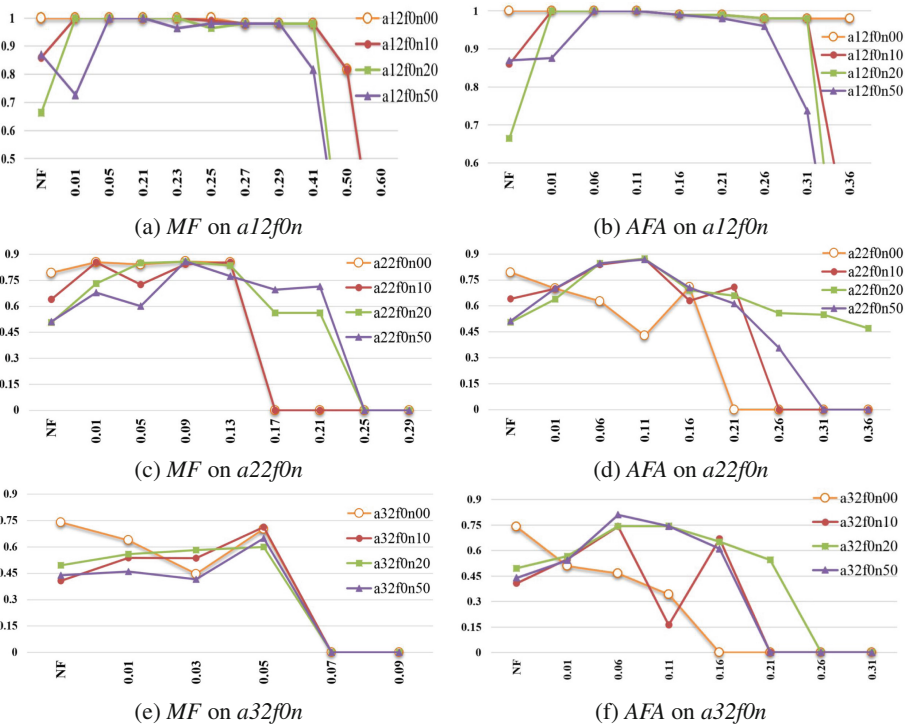
In Fig. 5, we compare the average required time of applying the process discovery algorithms with/without filtering methods. In this figure, the *y*-axis represents the time in milliseconds with logarithmic scale. According to this figure, filtering methods reduce the required time for discovering process models, because there are fewer traces in the filtered event logs. Although, in *AF* the discovery time reduction is higher, the filtering time for this method is much higher than *MF* method. Therefore, generally *MF* is faster than *AF*.

In the last experiment, to evaluate the ability of our proposed filtering method in detecting traces that contain outliers and corresponding effects on quality of



**Fig. 5.** Average of required time for process discovery with/without filtering.

process discovery algorithms, we use three synthetic event logs; *a12f0n*, *a22f0n* and *a32f0n*. These event logs are artificially added by 0, 10, 20 and 50 percent of different types of noise [2]. The last two characters of event log indicate the percentage of noise added to it, for example, *a22f0n20* correspond to *a22f0n* that contains 20% noise. Here, noisy event logs are used for process discovery and the original synthetic event logs (which contain no noise) use for computing F-Measure. Similar to the experiment in Fig. 3, *IMi* algorithm with 11 various



**Fig. 6.** Effect of filtering thresholds on F-Measures of synthetic event logs. *y*-axes are indicating values of best F-Measure and *x*-axes are showing the filtering thresholds.

internal noise thresholds has been used, but we show results of best F-Measure. The results of this experiment is presented in Fig. 6. According to this figure, F-Measures of models improve when applying filtering methods. This improvement is much more substantial for event logs that contain more percentage of noise.

For *a12f0n* event log which has the simplest structure among these event logs, both methods lead to similar results. However, for *a22f0n*, applying *MF* results in better F-Measures. Finally, in *a32f0n* that corresponds to the most complex model with lots of parallelism, *AFA* performs better than *MF*. This can be explained by the fact that when a lot of parallelism is present, the conditional probability of non-outlier behaviour is low as well, i.e. parallelism implies a lot of variety in behaviour. In this situations it seems using short subsequences (e.g.  $k = 1$ ) or applying smaller  $\kappa$  value are better choices for *MF*.

These experiments indicate that the proposed filtering method is useful for process discovery algorithms to have models with higher F-Measure and it reduces their required discovery time. This way tends to outperform state-of-the-art process mining filtering techniques.

## 6 Conclusion

Process discovery is used to extract process models from event logs. However, real event logs contain noise and infrequent behaviour that make to discovery process model from the whole event log problematic. Separating these outliers from event logs is beneficial for process discovery techniques and helps to improve process discovery results.

To address this problem, we propose a filtering method that takes an event log as an input and returns a filtered event log based on a given threshold. It usesn the conditional probability of occurrence of an activity after sequence of activities. If this probability is lower than the given threshold, the activity is considered as an outlier.

To evaluate the proposed filtering method we developed a plugin in the *ProM* platform and also offer it through *RapidProM*. As presented, we have applied this method on real event logs, and several process discovery algorithms. Additionally, we use the proposed method on three synthetic event logs. The results indicate that the proposed approach is able to help process discovery algorithms to discover models that better balance between different behavioural quality measures. Furthermore, using these experiments we show that our filtering method outperforms related state-of-the-art process mining filtering techniques.

We plan to evaluate the effect of using different values of  $k$ , i.e. length of subsequences. Also, other metrics besides F-Measure like simplicity, generalization and structuredness could be analyzed. we want to apply event-level filtering and also assess different ways of using  $\kappa$ . The other approach in this domain will be providing event-level filtering that we ignore it in this paper.

## References

1. van der Aalst, W.M.P.: *Process Mining - Data Science in Action*, 2nd edn. Springer, Heidelberg (2016)
2. Maruster, L., Weijters, A.J.M.M., van der Aalst, W.M.P., van den Bosch, A.: A rule-based approach for process discovery: dealing with noise and imbalance in process logs. *Data Min. Knowl. Discov.* **13**(1), 67–87 (2006)
3. Van der Aalst, W.M., van Dongen, B.F., Günther, C.W., Rozinat, A., Verbeek, E., Weijters, T.: Prom: The process mining toolkit. *BPM (Demos)* 489(31) (2009)
4. van der Aalst, W.M.P., Bolt, A., van Zelst, S.J.: RapidProM: Mine your processes and not just your data. *CoRR* abs/1703.03740 (2017)
5. van Zelst, S., van Dongen, B., van der Aalst, W., Verbeek, H.: Discovering Relaxed Sound Workflow Nets using Integer Linear Programming. *arXiv preprint arXiv:1703.06733* (2017)
6. Buijs, J.C.A.M., van Dongen, B.F., van der Aalst, W.M.P.: On the role of fitness, precision, generalization and simplicity in process discovery. In: Meersman, R., Panetto, H., Dillon, T., Rinderle-Ma, S., Dadam, P., Zhou, X., Pearson, S., Ferscha, A., Bergamaschi, S., Cruz, I.F. (eds.) *OTM 2012. LNCS*, vol. 7565, pp. 305–322. Springer, Heidelberg (2012). [https://doi.org/10.1007/978-3-642-33606-5\\_19](https://doi.org/10.1007/978-3-642-33606-5_19)
7. van der Aalst, W., Weijters, T., Maruster, L.: Workflow mining: discovering process models from event logs. *IEEE Trans. Knowl. Data Eng.* **16**(9), 1128–1142 (2004)
8. Weijters, A.J.M.M., Ribeiro, J.T.S.: Flexible Heuristics Miner (FHM). In: *IEEE Symposium on Computational Intelligence and Data Mining (CIDM)*. IEEE (2011)
9. van der Aalst, W.M.P., Rubin, V., Verbeek, H.M.W., van Dongen, B.F., Kindler, E., Günther, C.W.: Process mining: a two-step approach to balance between underfitting and overfitting. *Softw. Syst. Model.* **9**(1), 87–111 (2008)
10. Günther, C.W., van der Aalst, W.M.P.: Fuzzy mining – adaptive process simplification based on multi-perspective metrics. In: Alonso, G., Dadam, P., Rosemann, M. (eds.) *BPM 2007. LNCS*, vol. 4714, pp. 328–343. Springer, Heidelberg (2007). [https://doi.org/10.1007/978-3-540-75183-0\\_24](https://doi.org/10.1007/978-3-540-75183-0_24)
11. van der Werf, J.M.E.M., Dongen van Dongen, B.F., Hurkens, C.A.J., Serebrenik, A.: Process discovery using integer linear programming. *Fundam. Inform.* **94**(3–4), 387–412 (2009)
12. Leemans, S.J.J., Fahland, D., van der Aalst, W.M.P.: Discovering block-structured process models from event logs - a constructive approach. In: Colom, J.-M., Desel, J. (eds.) *PETRI NETS 2013. LNCS*, vol. 7927, pp. 311–329. Springer, Heidelberg (2013). [https://doi.org/10.1007/978-3-642-38697-8\\_17](https://doi.org/10.1007/978-3-642-38697-8_17)
13. Leemans, S.J.J., Fahland, D., van der Aalst, W.M.P.: Discovering block-structured process models from event logs containing infrequent behaviour. In: Lohmann, N., Song, M., Wohed, P. (eds.) *BPM 2013. LNBIP*, vol. 171, pp. 66–78. Springer, Cham (2014). [https://doi.org/10.1007/978-3-319-06257-0\\_6](https://doi.org/10.1007/978-3-319-06257-0_6)
14. van Zelst, S.J., van Dongen, B.F., van der Aalst, W.M.P.: Avoiding over-fitting in ILP-based process discovery. In: Motahari-Nezhad, H.R., Recker, J., Weidlich, M. (eds.) *BPM 2015. LNCS*, vol. 9253, pp. 163–171. Springer, Cham (2015). [https://doi.org/10.1007/978-3-319-23063-4\\_10](https://doi.org/10.1007/978-3-319-23063-4_10)
15. Yang, W., Hwang, S.: A process-mining framework for the detection of healthcare fraud and abuse. *Expert Syst. Appl.* **31**(1), 56–68 (2006)

16. Ghionna, L., Greco, G., Guzzo, A., Pontieri, L.: Outlier detection techniques for process mining applications. In: An, A., Matwin, S., Raś, Z.W., Ślęzak, D. (eds.) ISMIS 2008. LNCS (LNAI), vol. 4994, pp. 150–159. Springer, Heidelberg (2008). [https://doi.org/10.1007/978-3-540-68123-6\\_17](https://doi.org/10.1007/978-3-540-68123-6_17)
17. Wang, J., Song, S., Lin, X., Zhu, X., Pei, J.: Cleaning structured event logs: a graph repair approach. In: IEEE 31st International Conference on Data Engineering, ICDE, pp. 30–41 (2015)
18. Cheng, H.J., Kumar, A.: Process mining on noisy logs –can log sanitization help to improve performance? *Decis. Support Syst.* **79**, 138–149 (2015)
19. Cendrowska, J.: PRISM: An algorithm for inducing modular rules. *Int. J. Man Mach. Stud.* **27**(4), 349–370 (1987)
20. Conforti, R., La Rosa, M., ter Hofstede, A.H.M.: Filtering out infrequent behavior from business process event logs. *IEEE Trans. Knowl. Data Eng.* **29**(2), 300–314 (2017)
21. Bolt, A., de Leoni, M., van der Aalst, W.M.P.: Scientific workflows for process mining: building blocks, scenarios, and implementation. *Int. J. Softw. Tools Technol. Transfer.* **18**(6), 607–628 (2016). <https://doi.org/10.1007/s10009-015-0399-5>
22. Weerdt, J.D., Backer, M.D., Vanthienen, J., Baesens, B.: A robust F-measure for evaluating discovered process models. In: Proceedings of the CIDM, pp. 148–155 (2011)