



HAL
open science

Stable Heuristic Miner: applying statistical stability to discover the common patient pathways from location event logs

Sina Namaki Araghi, Franck Fontanili, Elyes Lamine, Uche Okongwu,
Frederick Benaben

► To cite this version:

Sina Namaki Araghi, Franck Fontanili, Elyes Lamine, Uche Okongwu, Frederick Benaben. Stable Heuristic Miner: applying statistical stability to discover the common patient pathways from location event logs. *Intelligent Systems with Applications*, 2022, 14, pp.200071. 10.1016/j.iswa.2022.200071 . hal-03604368

HAL Id: hal-03604368

<https://imt-mines-albi.hal.science/hal-03604368v1>

Submitted on 25 Mar 2022

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution - NonCommercial - NoDerivatives 4.0 International License



Contents lists available at ScienceDirect

Intelligent Systems with Applications

journal homepage: www.elsevier.com/locate/iswa

Stable heuristic miner: Applying statistical stability to discover the common patient pathways from location event logs

Sina Namaki Araghi^{a,*}, Franck Fontanili^a, Elyes Lamine^a, Uche Okongwu^b,
Frederick Benaben^a

^a Industrial Engineering Center of IMT Mines Albi, Albi 81000, France

^b Department of Information, Operations and Management Sciences, Toulouse Business School, 20 Bd Lascrosses, Toulouse 31000, France



ARTICLE INFO

Article history:

Received 31 August 2021

Revised 10 February 2022

Accepted 4 March 2022

Available online 8 March 2022

Keywords:

Process mining

Statistical stability

Healthcare processes

Patient pathways

Indoor localization systems

ABSTRACT

Purpose: The classic heuristic miner algorithm has received lots of attention in the healthcare sector for discovering patient pathways. The extraction of these pathways provides more transparency about patient activities. The previous versions of this algorithm receive an event log and discover several process models by using manually adjustable thresholds. Then, the expert is left with the difficult task of deciding which discovered model can serve as the descriptive reference process model. Such a decision is completely arbitrary and it has been seen as a major structural issue in the literature of process mining. This paper tackles this problem by proposing a new process discovery algorithm to facilitate patient pathways diagnosis.

Approach: To address this scientific challenge, this paper proposes to consider the statistical stability phenomenon in an event log, and it introduces the stable heuristic miner algorithm as its contribution. To evaluate the applicability of the proposed algorithm, a case study has been presented to monitor patient pathways in a medical consultation platform.

Originality: Thanks to this algorithm, the value of thresholds will be *automatically calculated at the statistically stable limits*. Hence, instead of several models, only one process model will be discovered. To the best of our knowledge, applying the statistical stability phenomenon in the context of process mining to discover a reference process model from location event logs has not been addressed before.

Findings/Practical implications: The results enabled to remove the uncertainty to determine the threshold that represents the common patient pathways and consequently, leaving some room for potential diagnosis of the pathways.

© 2022 The Author(s). Published by Elsevier Ltd.

This is an open access article under the CC BY-NC-ND license (<http://creativecommons.org/licenses/by-nc-nd/4.0/>)

1. Introduction

This research work introduces its practicality by considering two challenges, the business and technical problems. The business problem presents the main objective of this research project. The technical problem explains the focus of this paper.

* Corresponding author.

E-mail addresses: sina.namakiaraghi@mines-albi.fr (S. Namaki Araghi), franck.fontanili@mines-albi.fr (F. Fontanili), elyes.lamine@mines-albi.fr (E. Lamine), u.okongwu@tbs-education.fr (U. Okongwu), frederick.benaben@mines-albi.fr (F. Benaben).

URL: <https://www.imt-mines-albi.fr> (S. Namaki Araghi), <https://www.imt-mines-albi.fr> (F. Fontanili), <https://www.imt-mines-albi.fr> (E. Lamine), <https://www.tbs-education.fr/> (U. Okongwu), <https://www.imt-mines-albi.fr> (F. Benaben)

1.1. The business problem: how to diagnose deviations in patient pathways by using a descriptive reference process model?

Process mining applications for monitoring patient pathways have been seen as a valid solution for such issues (Garcia et al., 2019; Thiede, Fuerstenau, & Barquet, 2018). Researchers in the field of process mining are trying to diagnose these operational problems by using the indoor location data of patients (Fernandez-Llatas, Lizondo, Monton, Benedi, & Traver, 2015; Martinez-Millana et al., 2019).

We argue that we need a *common pathways* as a reference model in order to diagnose the unexpected deviations in each patient pathway. In practice, this reference model is provided by doing interviews with many different domain experts.

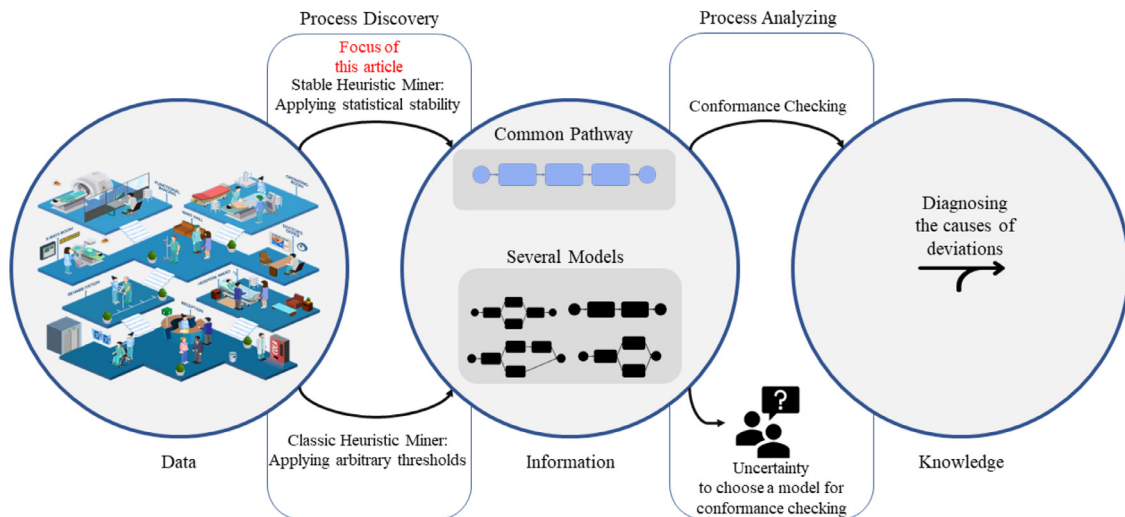


Fig. 1. Difference between the results of the heuristic miner algorithm and the stable heuristic miner algorithm presented in this paper.

This interview-based approach has its own shortcomings (Dumas, La Rosa, Mendling, Reijers et al., 2013).

Therefore, To diagnose patient pathways, we propose to automatically discover a descriptive reference process model. By comparing the descriptive reference process model with each patient pathway, we can detect possible deviations. Extraction of the descriptive reference process model is the main goal of this paper. Such a model represents the common pathways of patients in the hospital premises. However, there are several scientific obstacles that defy a process discovery algorithm to extract such a model.

1.2. The technical problem: how to discover the descriptive reference process model?

Van der Aalst presents a picture in van der Aalst (2016) to show the main challenges for the task of discovering a descriptive reference process model. Assuming that we are aiming to extract a so-called target model from an event log, we have to avoid certain outcomes. A “non-fitting” model will not be suitable, as it is not able to capture enough volume of information to express the main behavior patterns. Also, any process discovery technique should avoid “over-fitting” too. This means that an ideal model should not capture all of the behaviors in the event log. Additionally, the “under-fitting” concept is related to models which allow for the generation of behaviors that do not exist in the event log. Given these constraints, it is not clear how to find a trade-off among these notions. Under these circumstances, the author indicates that there is no exact definition of the target model and it is not clear how to respect all the criteria mentioned (van der Aalst, 2016). This is a valid issue in applications that associate ILS and process mining to extract the common pathways. In the literature of process mining Munoz-Gama et al. (2022), use of wearable devices (such as localization technologies) are identified as a mean to give different views of patients’ activities. These views can be unsteady and varying due to the exact technical problem defined above and shown in Fig. 1. This is a limiting factor for our research since we need to find a reference model that could be used for two purposes of diagnosing and simulating patients’ pathways.

For instance, consider the classic heuristic miner algorithm which is mostly used in the healthcare sector (Rojas, Munoz-Gama, Sepúlveda, & Capurro, 2016). This algorithm uses manually adjustable thresholds to discover the process models.

As shown in Fig. 1, these thresholds extract several process models with different levels of information. Therefore, it is de-

pendent on the experience of the domain expert to decide which model corresponds to the common pathways, or as called here the descriptive reference process model of patient pathways. Consequently, this uncertainty for selecting a reference model becomes a blocking point for further diagnostic actions. In addition, discovering a reference model could help us with simulation actions in future to better analyze and improve the processes. The manual adjustment of these thresholds has been seen as a structural challenge in the literature of process mining and it is an obstacle for the diagnostic and simulation actions (De Cnudde, Claes, & Poels, 2014; Janssenswillen, 2021). Fig. 1 summarizes the mentioned issues. With this in mind, the scientific question addressed in this paper is:

- How we can discover the common pathways¹ of patients without making arbitrary changes in the threshold values of the process discovery algorithm?

To answer this question, this paper introduces the Stable Heuristic Miner algorithm as its main contribution. This algorithm alters the classic heuristic miner algorithm by applying the criteria of the statistical stability phenomenon to discover the common pathways. Consequently, it would no longer be needed to determine the values of thresholds in a counter-intuitive way. It is important to mention why the application of statistical stability phenomenon seems to be a suitable solution. The reason is the fact that hospitals and in particular patient pathways are seen as examples of systems with emergent properties (Gorban, 2017). Therefore, in order to represent the common behavior of such systems, we need to consider the statistical stability phenomenon in the ensemble of patients processes (Gorban, 2014). We elaborate on this in Section 2.3.

This paper embraces the DIAG² research methodology (Namaki Araghi, 2019) to apply its method on the location data of patients. With an objective to diagnose and simulate business processes, DIAG methodology clarifies the necessary functions to transform raw location data into meaningful information to support decision making actions. Table 1 clarifies the definition of some of the used terms in this paper.

The remainder of this paper is structured as follows: in the second section—Material and methods—we will discuss similar studies

¹ The descriptive reference process model of patient pathways.

² Data state, Information State, Awareness, Governance.

Table 1
Definition of terms used in this paper.

Used term	Description
Relation frequency	Here it is defined as the number of times that the relation among activities is established. For instance, number of times activity 'c' is followed by 'd'.
Stable behavior	It represents behaviors that can be considered as the common behavior of the whole system with emergent properties. It is manifested by the statistical stability phenomenon.
Descriptive reference process model	It is used as the common pathways. It represents activities that are expressing the stable behaviors of patients while executing their processes.

on process discovery that are presented in the literature by previous authors. This discussion will be extended to cover the traditional approach of the heuristic miner algorithm, before introducing the statistical stability phenomenon. After demonstrating the scientific gaps, we will present a detailed illustration of the stable heuristic miner algorithm in the third section. Then, through the description of a case study in section four, the results of the proposed algorithm will be pragmatically evaluated and it will be discussed in section five. Finally, in the sixth and last section, we will conclude by presenting the strengths and limitations of this method as well as the future perspectives of this research work.

2. Material and methods

2.1. Related works

The literature of process discovery methods is extremely vast and rich. Authors in Augusto et al. (2018) provided a complete overview of the existing algorithms. Considering this technical problem presented at Section 1.2, many researchers proposed to find a process model that is easy to visualize Chapela-Campa, Mucientes, and Lama (2019); De San Pedro, Carmona, and Cortadella (2015); Vázquez-Barreiros, Mucientes, and Lama (2015). Even-though this is necessary for users to see a process model for understanding “what is happening”, this does not guarantee that the discovered model is showing the *reference model*.

Similar works in our area aim to extract a so-called *reference model* which represents the common behavior pattern recorded in event logs. Relevant discovery methods (for mining the reference model) in the literature are: model-based, clustering-based and profiling-based approaches (Li & van der Aalst, 2017). We will discuss these approaches in the following.

An example of model-based approaches is the work of Bezerra and Wainer where they propose the iterative and sampling algorithms for finding frequent cases and detecting variations. They identified the “dynamic threshold algorithm” for anomaly detection of traces in event logs (Bezerra & Wainer, 2013).

Clustering-based approaches are seen as more suitable than model-based ones for use in unstructured processes like patient pathways (Rebuge & Ferreira, 2012). However, these methods aim to detect clusters of behaviors rather than detecting the deviations that are causing the instability in the processes. Additionally, they can be time-consuming (Li & van der Aalst, 2017).

The model-based and clustering-based methods have been challenged in Li and van der Aalst (2017). Authors in Li and van der Aalst (2017) have indicated that these methods are either slow or inaccurate when dealing with complex event logs and unstructured processes that may contain *many activities*. Therefore, they proposed a novel profiling-based approach which creates a “profile” of cases that are representative of the majority of normal behaviors in the event log. Their approach has several steps: first, they sample all the cases in the event log. Then, based on a defined norm function they gather normal cases and identify them as the mainstream cases. Once they have found the mainstream, they compute the similarities between the mainstream and other cases. By creating the concept of a “profile” they classify cases with mutual fea-

tures. Then, their method quantifies the similarities based on the profile and identifies the normal cases and deviating cases. By adjusting the “norm function” one could increase (or decrease) the probability for sampling of the normal cases. This could be viewed as a disadvantage, since the decision to determine the value of the norm function could be completely arbitrary. This is the challenge we aim to overcome in this paper.

The Inductive Miner algorithm is one of the appreciated methods for discovering process models as well (Leemans, Fahland, & Aalst, van der, 2014). It tries to find the most prominent splits in event logs and it detects related operators to describe each split. However, the inductive miner method uses the concept of thresholds for activities and edges to propose different model. An example of this method is the inductive visual miner plugin in ProM tool³. Again, it is difficult to find which model could be the reference model for us.

Split Miner (Augusto et al., 2019) is another interesting research work with similar objectives to ours. Authors in Augusto et al. (2019) designed their approach based on a need to answer *complexity in discovered models* (so-called spaghetti-like models), *low-fitness* and *over-generalization*. They aim to extract a model that shows perfect fitness regarding both activities and edges. The similarity of their work and ours is in the final goal which is finding a reference model. However, their generally-applicable approach is to extract a model based on the capacities of activities and edges. On the other hand, we are trying a different path to find a reference model to represent the common pathways of patients. We address the nature of the healthcare organization which is a complex system with emergent properties (Aziz-Alaoui & Bertelle, 2009). Therefore, in order to represent what is the behavior of the whole system—the common pathways—we need to look for statistical stability in location event logs of patients. We will elaborate on this in Section 2.3 and develop our approach by considering the logic of heuristic miner algorithm as well.

2.2. Classic heuristic miner algorithm

The motive behind considering heuristic miner algorithm is its ability to deal with healthcare processes (Rojas et al., 2016). However, it is faced with many structural problems, such as arbitrary selection of the threshold values, dealing with variations and defining ways to filter noisy behaviors.

Previous authors in (De Cnudde et al., 2014; van der Aalst, 2016; Weijters, Maruster, & Aalst, 2004; Weijters & Ribeiro, 2011) have presented the heuristic miner algorithm by a series of steps to capture the behavior according to an event log (c.f. Fig. 2). The classic method contains five main steps in order to extract a process model which represents the behavior of an event log.

These steps are: (i) identify the footprint matrix, (ii) calculate the dependency measures, (iii) devise the graph, (iv) discover the splits and joins, and (v) adjust the mining loops with length of 1 and 2. This paper presents an alteration at the second step, which

³ <https://www.promtools.org/>.

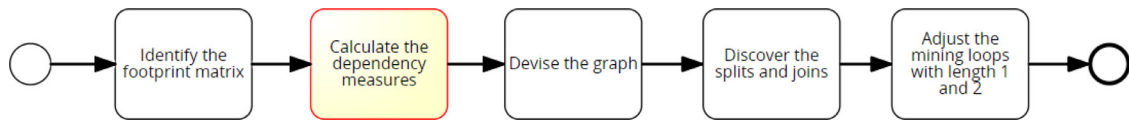


Fig. 2. Basic steps in the classic heuristic miner algorithm (van der Aalst, 2016).

is the core of heuristic miner algorithm. Addressing the last two steps is beyond the limit of this paper.

Authors in Weijters and Ribeiro (2011) define the dependency graph as the result of the first 3 steps. This term is identified as follows:

$$\text{Dependency Graph} = \{(a, b) \mid (a \in E \wedge b \in a\Box) \vee (b \in E \wedge a \in \Box b)\} \quad (1)$$

Here 'E' is defined as a limited set of activities. For this set of activities several events are recorded. ' $\Box b$ ' stands for the activities that come before 'b'. ' $a\Box$ ' denotes the activities that come after 'a'. Hence, in a dependency graph each activity can have input – output activities which are presented as a dependency relation (a, b). In order to devise the dependency graph, the number of times that an activity is directly followed by another one is presented in the format of a footprint matrix. Previously, heuristic miner algorithm aimed to define values among relationships known as “dependency measures”. These algorithms (Weijters & van der Aalst, 2003; Weijters & Ribeiro, 2011) would present different volumes of information in the process model by manually adjusting several thresholds within a range of –1 and 1. These values are calculated by this formula:

$$\text{Dependency Measure} : a \Rightarrow_w b = \frac{|a >_w b| - |b >_w a|}{|a >_w b| + |b >_w a| + 1} \quad (2)$$

where 'w' represents the event log with 'n' the number of activities and $|a >_w b|$ the number of times activity 'a' is followed by 'b'. The dependency measure value helps to determine the relationship between two certain activities. In the current application of these algorithms, experts use multiple thresholds. Since these thresholds are determined in an arbitrary way, the validity of the mined process model is dependent on the experience of its users. This is a major structural challenge for the classic heuristic miner algorithm (De Cnudde et al., 2014).

Therefore, in order to solve this issue, we propose to discover a state that represents the statistically stable behavior of the event log. The objective is to be able to discover a more structured process model and to give the algorithm greater ability to deal with complex processes and noise. This new approach removes the previous challenge of counter-intuitively choosing the value of thresholds. As a result, one can acquire a reference model for the further diagnostic actions. The new algorithm presented here has been inspired by the definition provided by Gorban where he declares that in order to represent the main behavior of a system with emergent properties such as a hospital, statistical stability needs to be found between the relation frequencies (Gorban, 2017).

2.3. Statistical stability phenomenon

It is best to illustrate this phenomenon by an example in nature. Consider a flock of birds in the sky, or a shoal of fishes in the sea. Their motions can be conceived as shapes. These shapes represent the behaviors of different groups which have emergent properties. The differences in these shapes are due to the different properties of the groups. Such behaviors (shapes) can be revealed by the statistical stability phenomenon (Gorban, 2017).

Important to consider that in each of these groups there are some existing deviating behaviors which at first are not seen by

the eye of an observer while looking at them from a distance, because these behaviors do not represent the stable behavior of their movements. Still, detecting these deviations is feasible.

The statistical stability of relation frequencies is an important property for analyzing the common behavior of a system with emergent properties such as hospitals. This phenomenon is manifested not only by considering the frequency of mass events, but also by the stability of the averages, the variances, and the standard deviations of the samples, and this is a feature that can be inherent in the collection of events (Gorban, 2017).

In essence, the definition of the statistical stability phenomenon can be inferred as:

In a system with emergent properties, there exists a state that shows a snapshot of the system behavior that could be seen as a common and stable representation –model– of the system in which the statistical stability is manifested Inspired by this phenomenon, this research work proposes a novel method based on the statistical stability phenomenon to discover the common pathways of patients from their movements in hospitals. This would help to remove the deviating behaviors which do not represent the common behavior of a group of patients and thus capture a descriptive reference process model. The next section will present how the statistical stability can be determined for patient pathways.

3. Theory

3.1. Preliminaries

Fig. 3 shows the steps in the new algorithm. Accordingly, an alteration is made by removing the manually configurable thresholds and replacing them by an action for statistically determining the thresholds from data.

This new modification would automatically detect—from relation frequencies in an event log—the activities that represent the statistically stable behavior while removing any unstable behaviors.

One of the methods used to demonstrate statistical stability is the creation of Shewhart control charts (Montgomery, 2007). These control charts must be devised mathematically and discovered from the event log. Eventually, they should indicate the thresholds of the statistical stable state.

Generally, control charts contain a center line that represents the average value of a measured characteristic, corresponding to the in-control state. Two other thresholds are called Upper Control Limit (UCL), and Lower Control Limit (LCL). These limits are calculated by considering the standard deviations and averages of the samples. These two limits (UCL, LCL) are the borders of a statistically stable state. As long as the graphed data points fall between these two thresholds the outcomes of the process are in-control. If a data point falls outside these limits, it will be considered as a variation of the process outcomes, and the process will no longer be considered stable. To apply this notion to discovering a stable process model, several assumptions and definitions have been considered.

3.1.1. Assumptions

Inspired by Shewhart control charts, three main assumptions have been made to find the statistically stable behaviors of processes from an event log.

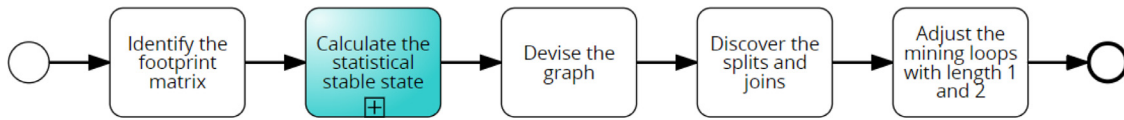


Fig. 3. The sequence of actions for applying stable heuristic miner.

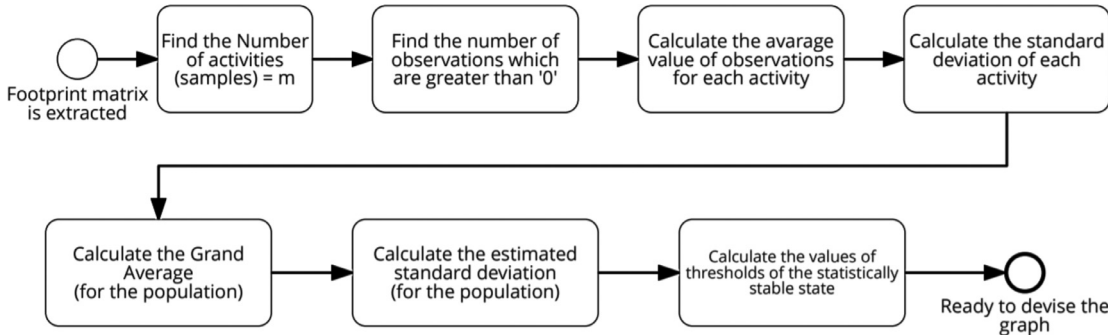


Fig. 4. The sequence of applying calculations for the stable heuristic miner algorithm.

1. The first assumption is the normality of the data distribution regarding the relation frequencies among registered activities in the event log. Most of the statistical methods for investigating the data are highly dependent on the distribution of data. If the data is not normally distributed, then certain adjustments in the method should be considered so as to adapt it to the distribution function of the data. This normality assumption is made due to the degree of freedom in which that data distribution could change. This assumption is authentic and justifiable by the Central Limit Theorem (Barany, Vu et al., 2007) and the statistical process control paradigm (Montgomery, 2007).
2. The second assumption is the one-sided-stability assumption. This means considering and presenting the activities that have an average of their relation frequencies greater than UCL. Usually, by applying the new method, activities with a high level of variations will be outside of UCL, which is mathematically correct. These activities are considered as not coming within the stable behavior of the log. However, such activities could provide information regarding the behaviors that show higher levels of variations and are thus the cause of the instability in the whole behavior. In order to not ignore these activities, they are displayed as “hot zones”, color-coded red. The reason behind this is that it is extremely rare to extract a pattern from an event log that shows all the activities illustrating a stable behavior. Therefore, the deviating activities that are causing instability in a model by a higher level of variations are presented in the process model.
3. In the footprint matrix the last activity will have '0' values, since it will not be followed by any other activity. Here, in order to not avoid the ending zone of the process, the last activity will be considered as an “end activity” with one “observation”.

3.1.2. Definitions and the sequence of functions in the algorithm

Fig. 4 illustrates each step of the algorithm. In the following we will elaborate on each step by using an example which is presented below as 'L'.

$L = [< a, b, c, d, e, l, m >^{12}, < a, b, f, d, e, l, m >^2, < a, b, c, d, g, e, l, m >$
 $< a, b, c, d, g, h, e, l, m >^5, < a, c, b, c, d, l, m >^6, < a, c, f, c, d, l, m >^3$
 $< a, c, b, i, c, e, c, d, g, l, m >^5, < a, c, b, c, f, c, d, l, m >^6, <$
 $a, c, b, c, i, c, h, c, d, l, m >^4, < a, c, b, c, i, f, c, d, g, l, m >^6, <$
 $a, b, c, j, l, m >^6, < a, c, b, k, e, d, l, m >^4]$

Within this example (eventlog 'L'), each group represents a trace. Every trace consists of events corresponding to the activities. For example, the first trace $< a, b, c, d, e, l, m >^{12}$ shows that 12 cases have followed the same sequence of activities. As shown in Fig. 4, the algorithm needs a footprint matrix as an input to extract and record the number of times one activity is directly followed by another activity. This matrix was presented previously in van der Aalst (2016); Weijters and Ribeiro (2011). The description of a footprint matrix is detailed within the definition 1.

Definition 1 (Population S). The footprint matrix here is considered as the “population” in which all of the relation frequencies are presented. This matrix represents the direct relations among different activities within an event log.

The matrix below shows the direct relations among all the activities recorded in the eventlog 'L'. For example, this matrix shows that activity 'a' is followed 26 times by activity 'b'. This matrix will be used as an input to initialize the algorithm.

$$S = \begin{pmatrix} & a & b & c & d & e & f & g & h & i & j & k & l & m \\ a & 0 & 26 & 34 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ b & 0 & 0 & 46 & 0 & 0 & 2 & 0 & 0 & 5 & 0 & 4 & 0 & 0 \\ c & 0 & 31 & 0 & 48 & 5 & 9 & 0 & 4 & 10 & 6 & 0 & 0 & 0 \\ d & 0 & 0 & 0 & 0 & 14 & 0 & 17 & 0 & 0 & 0 & 0 & 23 & 0 \\ e & 0 & 0 & 5 & 4 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 20 & 0 \\ f & 0 & 0 & 15 & 2 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ g & 0 & 0 & 0 & 0 & 1 & 0 & 0 & 5 & 0 & 0 & 0 & 11 & 0 \\ h & 0 & 0 & 4 & 0 & 5 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ i & 0 & 0 & 9 & 0 & 0 & 6 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ j & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 6 & 0 \\ k & 0 & 0 & 0 & 0 & 4 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ l & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 60 \\ m & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{pmatrix}$$

Step1: Find the number of activities (samples)

According to Fig. 4, at first, we need to extract the number of samples (activities).

Definition 2 (Sample s). Each row ($i = a \rightarrow m$) in the population presents a vector that shows the relation frequency of an activity being followed by another one. Therefore, the corresponding vectors for the activities in the event log are considered as the samples of the population.

Table 2
The calculation of \bar{x} , C_{4n} , and σ for each sample and the \bar{x} of the population.

Activities	(\bar{x}_i)	(σ)	C_{4n}
a	$\bar{x}_1 = \frac{26+34}{2} = 30$	5.65	0.80
b	$\bar{x}_2 = \frac{46+2+5+4}{4} = 14.25$	21.2	0.92
c	$\bar{x}_3 = \frac{31+48+5+9+4+10+6}{7} = 16.14$	16.82	0.96
d	$\bar{x}_4 = \frac{14+17+23}{3} = 18$	4.5	0.88
e	$\bar{x}_5 = \frac{5+4+20}{3} = 9.6$	8.9	0.88
f	$\bar{x}_6 = \frac{15+2}{2} = 8.5$	9.19	0.80
g	$\bar{x}_7 = \frac{1+5+11}{3} = 5.6$	5.03	0.88
h	$\bar{x}_8 = \frac{4+5}{2} = 4.5$	0.7	0.80
i	$\bar{x}_9 = \frac{9+6}{2} = 7.5$	2.12	0.80
j	$\bar{x}_{10} = \frac{6}{1} = 6$	0	0
k	$\bar{x}_{11} = \frac{4}{1} = 4$	0	0
l	$\bar{x}_{12} = \frac{60}{1} = 60$	0	0
m	$\bar{x}_{13} = \frac{0}{1} = 0$	0	0

Grand average ($\bar{\bar{x}}$) = 14.22.

For example, in the matrix above there are 13 samples (vectors related to each activity = s, which relates to $i = a$ to m). This number is identified by 'm'. In this example $m = 13$.

Step2: Find the number of observations which are greater than 0

Definition 3 (Observations). The values within the population (footprint matrix) are considered as the observations. Each observation is identified by ' x_{ij} ', where 'i' stands for the rows in the footprint matrix and 'j' represents the columns.

These observations present the relation frequencies among existing samples. The total number of observations within the population is identified by 'N'. For the example shown here, $N = 30 + 1$. The '+1' is an adjustment that has been made based on the third assumption, that the "end activity" is not considered as a null sample.

Step3: Calculate the average value of observations for each activity

As the title of this step indicates, we simply measure the average value of observations for each sample (activity). The results are represented within **Table 2** under column (\bar{x}_i).

Step4: Calculate the standard deviation value of observations for each activity

As shown in **Table 2**, we indicate the value of standard deviation (σ) for each activity.

Eq. (3) shows the basic method for calculating the standard deviation of each sample. Note that x_{ij} is an observation within a sample; n_i is the size of the sample, and \bar{x}_i is the average of observations for the i th sample.

$$\sigma_i = \sqrt{\frac{(1)}{n_i - 1} \sum_{j=1}^m (x_{ij} - \bar{x}_i)^2} \tag{3}$$

Step5: Calculate the grand average

It is very important to note that here the sizes of the samples are not necessarily similar. For example, the sample size for activity 'a' is equal to 2 ($n_{s_a} = 2$) and for activity 'b' is ($n_{s_b} = 4$). Therefore, to make sure about unbiased factor of our analysis we need to need to define two important definitions. These definitions are the grand average, and the estimated standard deviation.

Definition 4 (Grand average. $\bar{\bar{x}}$) Since the size of samples is a changing variable, $\bar{\bar{x}}$ has been defined to express the average of relation frequencies (x_{ij}) within the whole population. **Eq. (4)** shows how the grand average would be calculated. Note that m is the

number of samples and \bar{x}_i stands for the average of the 'ith' sample.

$$\bar{\bar{x}} = \frac{\sum_{i=1}^m n_i \bar{x}_i}{\sum_{i=1}^m n_i} \tag{4}$$

For this example, the grand average of the entire population is equal to '14.22'.

Step6: Calculate the estimated standard deviation

Definition 5 (Estimated Standard deviation. $\hat{\sigma}$) Similarly to the grand average, the estimated standard deviation $\hat{\sigma}$ has been defined by **Eq. (6)** in order to understand how the behavior within the population deviates from one sample to another while the sample sizes differ.

To calculate the estimated standard deviation ($\hat{\sigma}$), we need to ensure about the unbiased calculation. Therefore, use a factor known as C_{4n} (**Montgomery, 2007**). This factor is dependent on the size of each sample.

In order to measure the C_{4n} factor for each sample, **Eq. (5)** will be used.

$$C_{4n} = \frac{4(n - 1)}{4n - 3} \tag{5}$$

For the mentioned example, the values of C_{4n_i} and σ_i are presented in **Table 2**.

Now by considering these values, the **Eq. (6)** is used to calculate the estimated standard deviation of the population:

$$\hat{\sigma} = \frac{1}{m} \sum_{i=1}^m \frac{\sigma_i}{C_{4n_i}} \tag{6}$$

In this example, the value of $\hat{\sigma}$ is equal to '6.42'.

$$\hat{\sigma} = \frac{1}{13} \times \left[\frac{5.65}{0.79} + \frac{21.20}{0.92} + \frac{16.82}{0.95} + \frac{4.5}{0.88} + \frac{8.96}{0.79} + \frac{9.19}{0.79} + \frac{5.03}{0.88} + \frac{0.70}{0.79} + \frac{2.12}{0.797} + 0 + 0 + 0 + 0 \right] = 6.42 \tag{7}$$

After acquiring these metrics, the algorithm can construct the control limits (thresholds) required to extract the stable behavior.

Step7: Calculate the values of thresholds of the statistically stable state

Definition 6 (Central Line(CL)). As it's shown in **Eq. (8)**, 'CL' represents the most stable behaviors. The activities whose average of recorded observations is close to CL are normally present in most traces of the log, and they are thus the core activities.

$$CL = \bar{\bar{x}} \tag{8}$$

In this example, CL is equal to 14.22.

On the other hand, in the literature on the subject of determining the stable state, a certain distance from the CL is allowed (**Montgomery, 2007**). Previously, a distance of 3σ was used for samples with unique sizes. In this paper, since the sample sizes are changing variables, the distance is defined by considering the population estimated standard deviation and two defined constants ($A_{3\bar{n}}C_{4\bar{n}}\hat{\sigma}$).

This distance helps to define the two other limits or borders of the stable state.

Definition 7 (Lower Control Limit (LCL)). This threshold filters the behaviors that do not represent the main and stable behavior of the log. The activities with an average (\bar{x}_i) lower than LCL will not be shown in the process model.

Therefore, the algorithm will consider the activities that have a stronger presence in the behavior within the event log and will

remove the deviations. Eq. (9) will use the previous definitions to determine this threshold.

$$LCL = \bar{x} - (A_{3\bar{n}} \times C_{4\bar{n}} \hat{\sigma}) \quad (9)$$

Note that $A_{3\bar{n}}$ is a customary constant for considering a previously defined distance from CL (Montgomery, 2007). It can be calculated by Eq. (10).

$$A_{3\bar{n}} = \frac{\bar{n}}{C_{4\bar{n}} \sqrt{\bar{n}}} \quad (10)$$

Since the thresholds apply to the whole population, the formula would consider the average of all the sample sizes for calculating the $A_{3\bar{n}}$ and $C_{4\bar{n}}$ factors. Therefore: $\bar{n} = \text{Average of sample sizes} = \frac{N}{m}$ which is equal to the total number of observations divided by the total number of samples.

In this example, the value of m is equal to 13. Also, the number of observations has been indicated: $N = 31$. Therefore, $\bar{n} = \frac{31}{13} \approx 3$. And, $A_{3\bar{n}} = 1.94$.

Definition 8 (Upper Control Limit (UCL)). As shown in Eq. (11), the value of UCL sets the bar for activities with the maximum amount of variations with regard to the whole population.

$$UCL = \bar{x} + (A_{3\bar{n}} \times C_{4\bar{n}} \hat{\sigma}) \quad (11)$$

As a result, activities with \bar{x} greater than UCL are considered here as the zones where their behavior causes the process to be unstable. This could lead to bottlenecks at these activities while executing the process. Such activities would generate behaviors that do not normally correspond to the behavior of the whole population.

Consequently, in this example, these thresholds are equal to:

- $UCL = 14.22 + (1.94) \times (0.88) \times (6.42) \approx 26$
- $CL = 14.22$
- $LCL = 14.22 - (1.94) \times (0.88) \times (6.42) \approx 4$

Definition 9 (Statistically stable state). Finally, Eq. (12) determines which activities express a stable behavior in accordance with the whole recorded information in an event log.

$$LCL < \bar{x}_i < UCL \quad (12)$$

Normally, if all of the recorded activities in an event log express a stable behavior, no activity will be removed. This could imply that the process is running smoothly. But, if a variation exists in the behaviors, it will be detected by means of the two thresholds (UCL, LCL).

Definition 10 determines which activities will be considered within the modeled common behaviors of the event log.

Definition 10 (Descriptive reference process model \mathcal{P}). The descriptive reference process model or the ‘‘common behaviors’’ will contain activities that respect the following conditions:

$$\begin{aligned} & [\forall \mathcal{A} \exists s] \wedge [\forall s \subseteq S \exists \bar{x}_s] \\ \therefore & (\mathcal{A} \in \mathcal{P}) \rightarrow [LCL < \bar{x}_s < UCL] \cup [UCL \leq \bar{x}_s] \end{aligned} \quad (13)$$

Definition 10 states that for each activity (\mathcal{A}), a vector of relation frequencies with other activities exists. This vector is defined as a sample of the population (footprint matrix). And, for each sample there exists a \bar{x}_s which represents the average of relation frequencies. Therefore, the corresponding activity to the sample (s) will be represented in the descriptive reference process model (\mathcal{P}) if the average of its relation frequencies is between the two thresholds (considered as stable behavior) or if it is greater than the UCL value (considered as the hot zones).

The steps of the stable heuristic miner algorithm could be realized by the set of algorithms presented in the (Appendix A: ‘‘extract the thresholds’’ and Appendix B: ‘‘identify the status of activities’’). In the next section, we will present the final outcome.

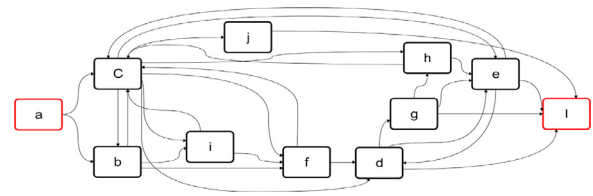


Fig. 5. An illustrative example of the algorithm outcome.

3.2. Result of the illustrative example

Concerning the example in Section 3.1.2, Fig. 5 illustrates the descriptive reference process model for the example log (L). This model represents the stable behavior of the example event log (L). Red activities (a and l) correspond to high variation in behaviors. Two activities (m and k) are removed with lower significance level for the general behavior.

4. Experimentation and results

4.1. The case study scenario

During this concrete example, 7 scenarios in 7 departments of a living lab of Toulouse hospital university were simulated. The selection of those departments and patient profiles was mainly based on the possibility and acceptance of the hospital in giving us the privilege to use its facilities and resources. The Toulouse Hospital University is located in south of France with several establishments. More than 3900 physicians and 11,600 hospital staff are welcoming around 280,000 patients annually. It has been estimated that more than 800,000 medical appointments are being registered each year. Approximately, 400 patients are admitted to the emergency department daily.

This experiment resulted in the generation of location data from the simulated healthcare processes of 261 patients. The stable heuristic miner algorithm is developed within an application known as R.IO-DIAG,⁴ to visualize the results of the experiment.

There were two objectives for this experiment: first, to obtain a descriptive reference process model showing the normal and stable pathways for all the departments in the hospital. This allows to visualize the zones that are being occupied by patients during the execution of healthcare processes. Secondly, experts tried to evaluate the descriptive reference process model for each department.

Several steps were taken prior to begin the experiment. At first, the primary information such as the maps, resources, zones information, and required patients' information were gathered. Then, this information was imported into the localization system. Beginning the experiments at this point would have led to a set of primary event logs which are not easily understandable, as they would only contain the location data of objects (x, y, z). Therefore these event logs needed to be prepared before importing them into R.IO-DIAG by defining what the objects in the process were and which location data corresponded to which zone. In addition, the activity that could occur in a zone needed to be defined too. Therefore, a primary knowledge was given to the system regarding these needs. This configuration is explained in previous research works (Araghi et al., 2018).

Later on, each patient received a tag. Each tag had an identification number corresponding to the different patients. After gathering the location data, the event logs were interpreted by the R.IO-DIAG location data interpreter. Meanwhile, the event logs were re-

⁴ <https://research-gi.mines-albi.fr/display/RIOSUITE/R-IOsuite+Home>.

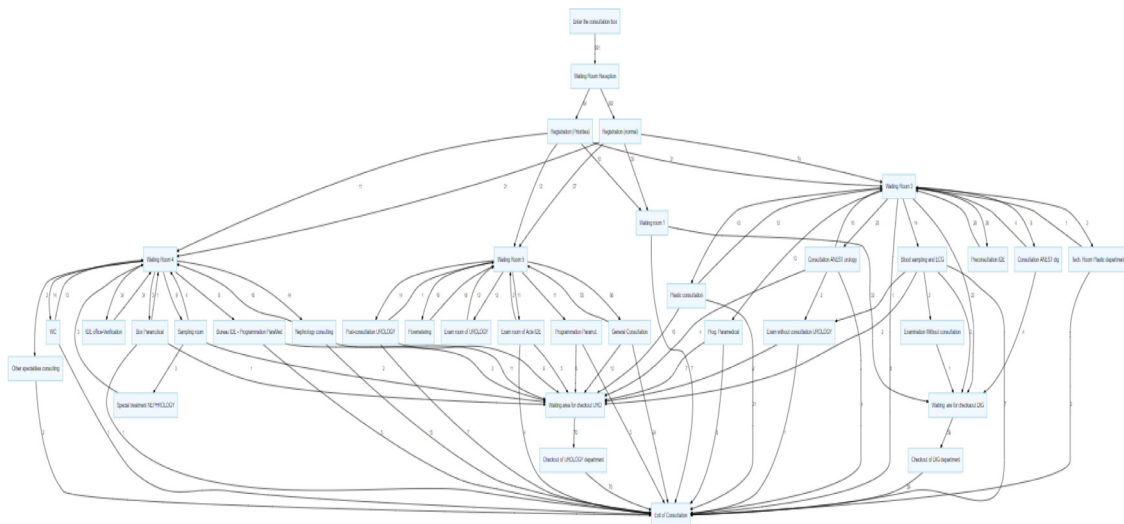


Fig. 6. A screen shot presenting patient pathways of all of the departments extracted by the classic heuristic miner algorithm.

fined. This action was necessary because several cases did not have complete data regarding their processes.

4.2. Process discovery results and discussion of the experiment

After this step, we asked the domain experts to analyze and diagnose the potential problems in patient pathways by using both the classic heuristic miner and the novel stable heuristic miner algorithms.

4.2.1. Analysis of all the departments

There were 36 activities registered in the main event log. At first, the experts used the classic heuristic miner algorithm. As presented in Section 2.2, this algorithm allows experts to extract several process models that each shows a different version of patient pathways. For instance, Fig. 6 is one of the outcomes of the classic heuristic miner algorithm by showing 100% of the registered information and it's not suitable nor visible for further analyses.

Other results of the classic heuristic miner algorithm for this case study can be seen online as supplementary materials (which are presented within this page⁵).

Each model represents different amount of information. Afterwards, we asked experts to highlight which model represents the common pathways of patients; so, we can use that model as a reference to diagnose potential problems and deviations in patient processes. As can be seen in those models (presented as supplementary online materials), each expert had different opinions about determining the common pathways. In line with what we have observed in the literature, we confirm that such a decision is completely dependent on the expert's experience. Moreover, the decision to select among these process models leads to uncertainty for the further diagnostic actions. We asked the experts to highlight (in the model they had chosen) any abnormal behavior that could lead to certain inefficiencies in the process; such as, an increase in waiting time, or poor resource allocation. The experts failed to reach a decision on these matters by using the results of the classic algorithm.

This observation led us to conclude that the classic heuristic miner helps to obtain a rapid illustration of patient pathways with

some flexibility to visualize the present activities in patient processes. However, our knowledge to diagnose the patient pathways is limited at this point.

Following these results, we proposed to the experts to use the stable heuristic miner algorithm. Consequently, the novel method represented one model that is extracted and evaluated by the logic of statistical stability phenomenon (c.f. Section 2.3). Fig. 7 shows this model, which represents the common pathways of patients. Thanks to the result of this algorithm, the experts observed which activities are normally present in patient pathways.

This model permits the experts to detect unstable activities and zones in the patient pathways, which was not possible by using the classic algorithm. As shown in Table 3, out of the total number of 36 activities in the event log, 16 are detected with an instability lower than the lower control limit (LCL) and are not shown in the descriptive reference process model (c.f. Fig. 7). From the 20 remaining activities in the descriptive reference process model, 7 are considered as hot zones, which impose high instability and variation to the normal behavior of the process.

These hot zones are indicated in red. These are the activities whose average behavior values are higher than the upper control limit (UCL). This implies that such activities represent unusual and eccentric behaviors in the log and this could lead to future problems.

To exemplify these statements, the "waiting_room_5" in Fig. 7 provides a good illustration. Based on the statistical stability, one can ensure that the probability of receiving the same behavior for this activity is high and all the activities related to "waiting_room_5" could be regenerated in the future. Therefore, the experts can plan and allocate the requirements for running such activities in the future. On the other hand, hot-zone activities such as "Registration_Normal" indicate that these activities are generating behaviors that are beyond the usual and stable behavior of all the other defined activities.

As a further illustration, the incoming flow into the "Registration_Normal" activity can be considered in Fig. 7. This activity is shown as a "hot zone" imposing high instability into the process. The reason is that 193 cases enter this activity, which is higher than most of the existing flows in the process. Also, this outgoing flow has a high variation in comparison with the outgoing flows from "Registration_Normal". Therefore, such behaviors could cause potential bottlenecks in this activity and consequent instability in the process. Similarly, the activity "Waiting_room_Reception" receives 261 cases and has two outgoing flows with values of 193

⁵ <https://research-gi.mines-albi.fr/display/gindresearch/Classic+Heuristic+Miner+results>

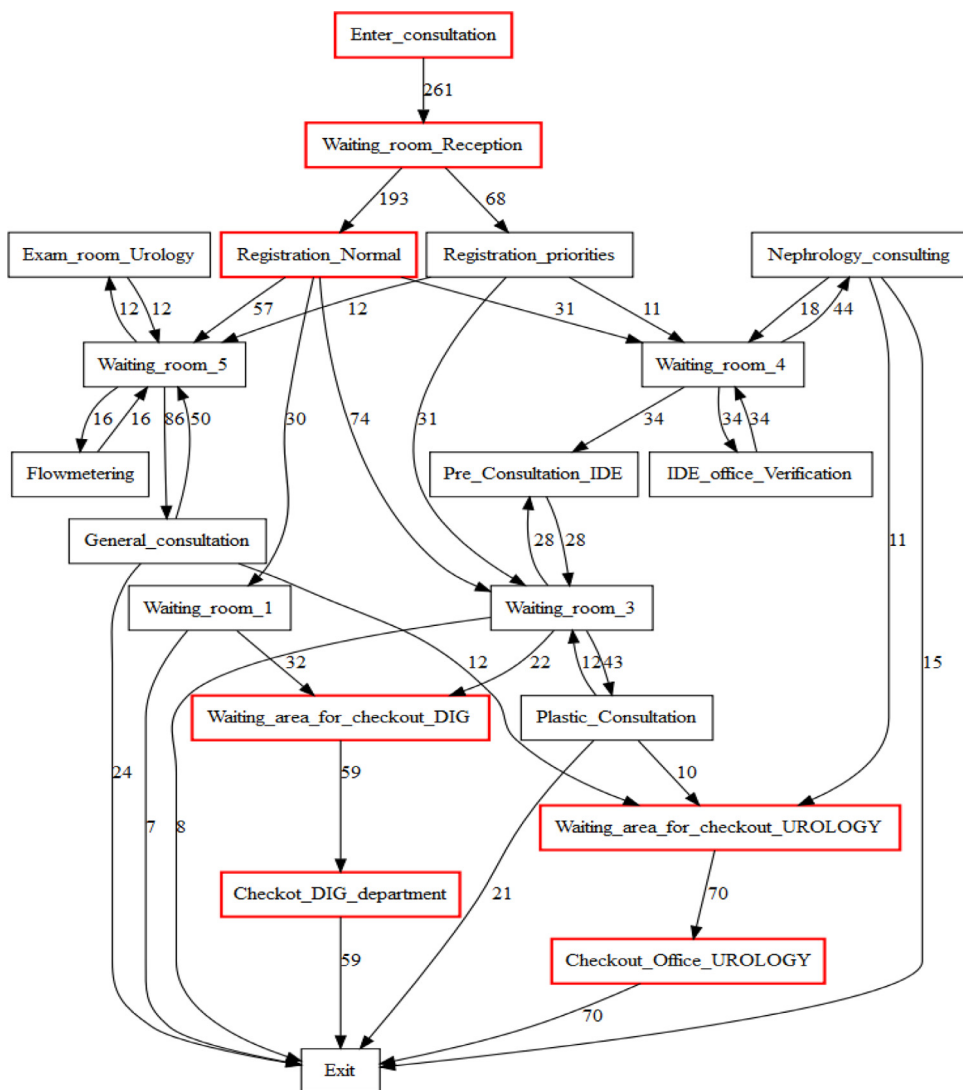


Fig. 7. A screen shot presenting the descriptive reference process model (common pathways) of all of the departments extracted by the stable heuristic miner algorithm.

Table 3

A comparison between the number of observed behaviors in the event log (of all the departments) and the modeled behaviors in the descriptive reference process model.

	Lower than LCL	In the stable state	Higher than UCL
Number of activities	16	13	7
Total number of modeled activities		20	
Total number of observed activities		36	

and 68. As shown in Fig. 7, the variation among these behaviors is significant. As a result, the unstable behavior of this activity is detected. These are the types of information and analyses that the stable heuristic miner algorithm promptly provides for the experts that the classic heuristic miner algorithm could not.

4.2.2. Analyzing one department

To address the second objective of this experiment, the processes of each department were investigated individually. As an example, the process model shown in Fig. 8 can be considered. It shows the patient pathways for the urology department according to the total existing events. This model is extracted by the classic heuristic miner. Similarly to the example that represented the total pathway of patients (c.f. Fig. 8), it is not clear which level of information represents the stable behavior.

In order to mine the descriptive reference process model for this department, the stable heuristic miner was used. Fig. 9 shows the common behavior of patients within the urology department of the hospital. From the 14 existing activities in the event log, 13 of them are shown within the descriptive reference process model of the urology department. Ten of these activities were detected as activities with stable behaviors and 3 of them (“Enter_consultation”, “Registration_Normal”, and “Exit”) show high instability in comparison with the total number of recorded behaviors.

In addition to these analyses, experts assumed that urology patients would carry out their administrative activities completely within the department. Despite this fact, one of the interesting results here is the absence of the checkout activity (“Checkout_Office_UROLOGY”) at the end of the descriptive reference pro-

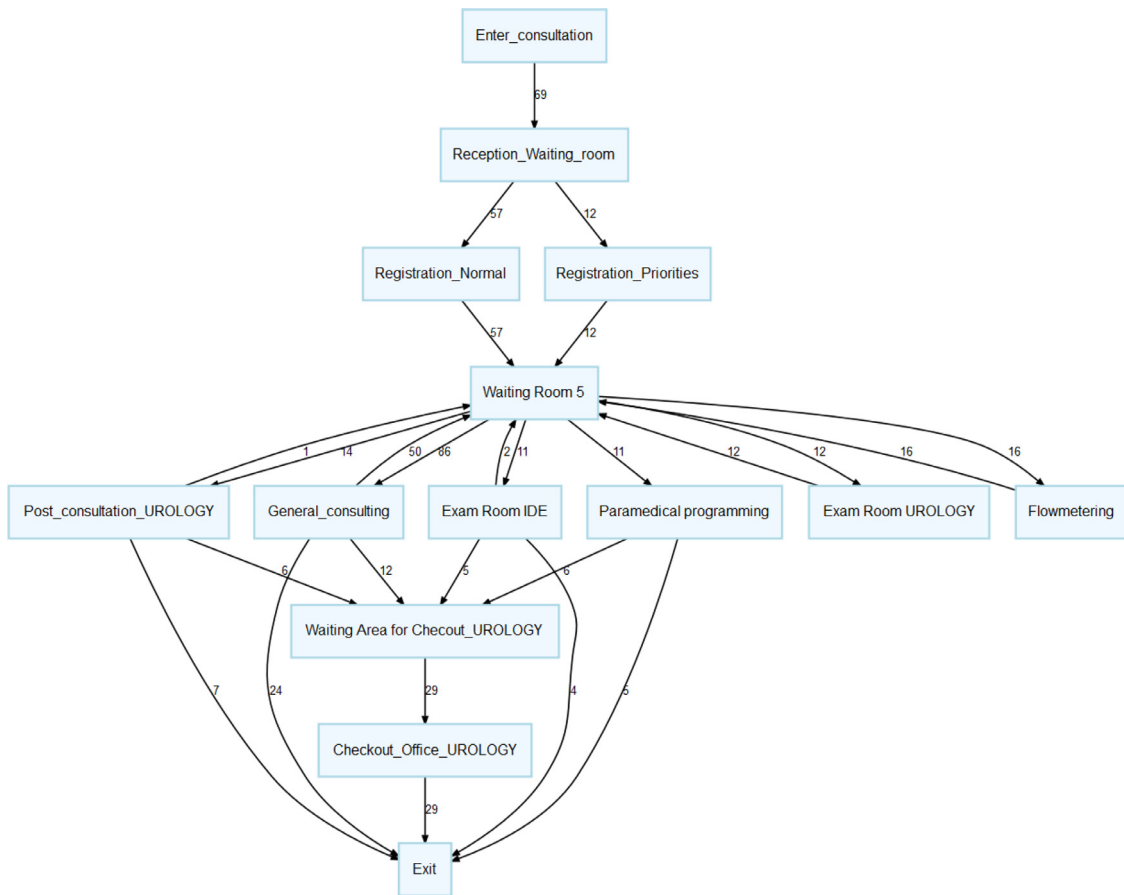


Fig. 8. The process model presenting patient pathways of the urology department discovered by the classic heuristic miner approach (thresholds set to show 100% of registered information).

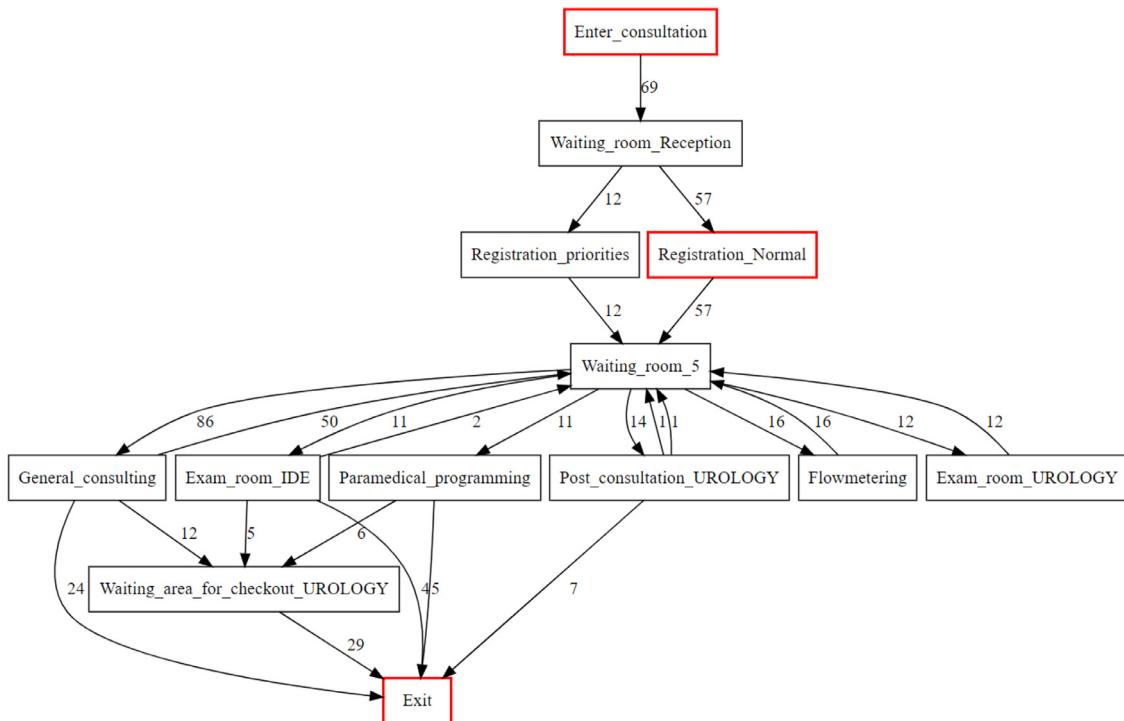


Fig. 9. A screen shot presenting the descriptive reference process model of the patients' pathway in the urology department extracted by the stable heuristic miner algorithm. Activities with high instability are indicated by a red color. (For interpretation of the references to color in this figure legend, the reader is referred to the web version of this article.)

Table 4

A comparison between the number of observed behaviors in the event log of the Urology department and the modeled behaviors in the descriptive reference process model.

	Lower than LCL	In the stable state	Higher than UCL
Number of activities	1	10	3
Total number of modeled activities		13	
Total number of observed activities		14	

cess model. The extracted model by the classic approach (c.f. Fig. 8) failed to detect such a shift in the process. But, by analyzing the discovered common pathway (c.f. Fig. 9), experts were able to highlight that patients did not often perform their checkout activities in this department. This was detected as an odd behavior in patient pathways. Such a deviation needed further analysis to find its causes. In this case, the deviation was due to the lack of resources in the other departments. Therefore, patients of other departments were being asked to perform their checkout in the urology department. This caused this zone to become vulnerable, leading to an unstable behavior. As a result, waiting time increased for this activity and consequently some patients avoided performing it altogether. Table 4 details the differences between the result of analyzing Urology event log by using both classic heuristic miner and the new stable heuristic miner algorithms.

5. Discussion

Acquiring such diagnoses is only feasible if experts are sure that the extracted model does indeed show the descriptive reference process model of patients within this department. This was not a possible outcome of previous algorithms that used location data of patients to discover the common pathways. In this paper, the application of the stable heuristic miner helped experts to detect the deviating behaviors automatically and to capture an image of what patients do normally, even if the experts did not particularly have complete knowledge of the process. Possessing such information allows the healthcare experts to avoid dissatisfying experience for patients. Moreover, it helps to detect patients who showed up for unscheduled activities.

There are evident advantages in discovering a descriptive reference process model Augusto et al. (2019); Estrada-Torres et al. (2021); Leno, Polyvyanyy, Dumas, La Rosa, and Maggi (2020). With an eye on diagnosing patient pathways, it is vital for the overall objective of our research to extract the common pathways of patients. Nonetheless, current methods did not satisfy our requirements. We are not only looking at visualizing the patient pathways, but we need to capture a reference model from location event logs so we can use them in diagnostic and simulation of patient pathways.

Indeed it is a difficult task to evaluate such an approach with quantified measures. Additionally, there is no common framework for evaluating such process discovery algorithms (De Cnudde et al., 2014; Estrada-Torres et al., 2021). Certain research works used conformance checking methods to evaluate the outcomes of the process discovery algorithms (Augusto et al., 2018). However, applying these methods (such as: *precision, generalization, simplicity, recall, fitness*) for evaluating stable heuristic miner imposed more ambiguity since they are not considering the stability among their evaluation criteria. It is a difficult task to find a trade off among all these criteria. In addition, the stable heuristic miner uses a method to evaluate statistical stability in an event log. Challenges regarding evaluation of process discovery algorithms have been seen in other research works as well. For instance in Estrada-Torres et al. (2021), the authors tried to extract a reference model so they can use it as a simulation model. They have applied Split Miner on multiple event logs in different domains so they can evaluate their ap-

proach. At the end, they have mentioned the need for an empirical evaluation since it is not evident how the current quantifiable criteria can evaluate the goodness of fit of a reference model. Due to such issues, we have used several event logs to discover the common patient pathways. Evaluating the new algorithm by using random available event logs is not a good practice in our case. Since, we are required to ensure that the system is an example of systems with emergent properties. This is a prerequisite to apply this method.

The presented analyses of the case study are due to the fact that we tried to manifest the basic definition of statistical stability phenomenon within the context of a location event log. Therefore, we looked at the event log as a population, and we considered each activity as a sample of the population. Then, we examined the statistical stability among relative frequency of events, average and standard deviations of direct relations between each activity. Fig. 10, presents a SWOT analysis (strength, weakness, opportunity, threat) Brender (2006) to elaborate on the results of the experiment while applying the new algorithm.

It is important to mention that the presented approach in this paper is not considered as a solution for an optimization problem. Previous works have applied meta-heuristic methods to address the challenges of process discovery algorithms (van der Aalst, de Medeiros, & Weijters, 2005). In case of our research, it is not possible to define the patient pathways discovery as an optimization problem. This is due to the nature of these processes. Patients can make different decisions to change their pathways and modify the objectives of a process.

6. Conclusion

In the context of monitoring and diagnosing patients processes, this paper sought a solution to overcome the challenge of automatically discovering a descriptive process model that could serve as the common patient pathways. To do so, the classic heuristic miner (Weijters & Ribeiro, 2011) was initially selected in this research work due to its abilities in providing satisfactory results for monitoring healthcare processes (Rojas et al., 2016). Subsequently, an obstacle emerged in selecting a model as a reference. Traditionally, this algorithm uses manually configurable thresholds for extracting different levels of information from an event log. Previously, this has been mentioned in the literature as an unsolved issue for this algorithm (De Cnudde et al., 2014). To address this issue, this paper applies the statistical stability phenomenon to evaluate the stability within all the existing relationships among activities. As a result, the new stable heuristic miner algorithm discovers the descriptive reference process model and removes the need to manually determine the thresholds.

We evaluated this algorithm by using an experiment. As mentioned in Sections 4.2.1 and 4.2.2, the classic approach resulted in the need to analyze several complicated models, and there was uncertainty about how to decide which model could be represented as the reference of common behaviors according to the total registered events. In contrast, the new stable heuristic miner algorithm directly provided the process which its structure was evaluated based on the statistical stability phenomenon. Moreover, if

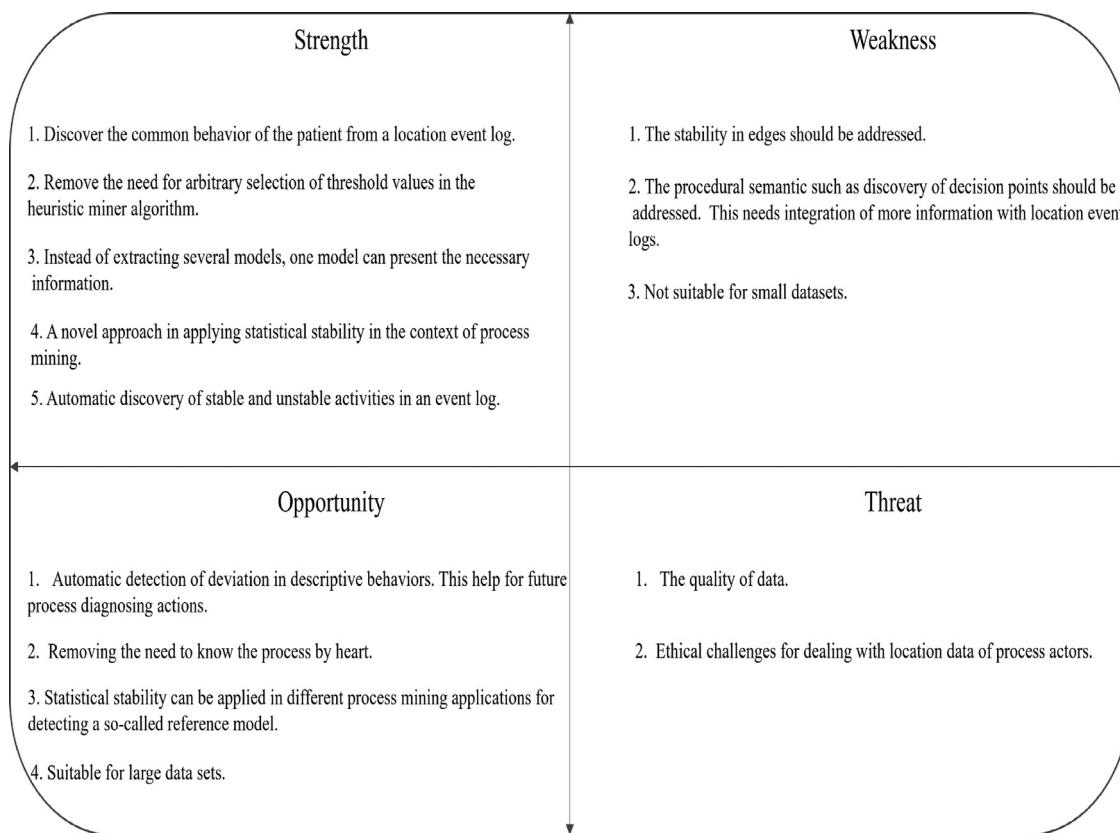


Fig. 10. A SWOT analysis of the proposed method in this paper.

the model was unstable, the algorithm revealed which activities were causing such instabilities.

6.1. General limitations

Wearable devices have shown many advantages in the healthcare sector (Corchado, Bajo, de Paz, & Tapia, 2008; Thibaud, Chi, Zhou, & Piramuthu, 2018; Zhou & Piramuthu, 2010). However, similar to most process mining applications, the proposed method here is also highly dependent on the quality, accuracy and reliability of the data. For instance, in a location event log, some cases may have disruptions in their data.

The stable heuristic miner algorithm needs to be improved so that it can discover the decision points (gateways) from the location data. This is a major challenge due to the lack of sufficient information in location event logs. It could be achieved by integrating other information from the hospital information system. Currently, this method considers activities representing the stable behavior of patients. However, it should also examine the statistical stability among edges (connections between activities). This limitation is visible in Fig. 5, where the edges illustrate a complex behavior. Additionally, it is important to evaluate this method outside the context of healthcare. However, as it was mentioned earlier in the introduction, the current method is addressing one challenge of a bigger research question in diagnosing patient pathways.

The true nature of statistical stability phenomenon is still somehow unclear. Nevertheless, lack of a clear comprehension of a

physical phenomenon is not a barrier for constructing theories to manifest them (Gorban, 2017).

6.2. Future perspective

It is important to devise a quantifiable evaluation method for discovering statistical stability in an event log. We have also made an assumption that data is normally distributed. The distribution of data can be at a constant change, which is valid for any dynamic systems with emergent properties. As a result, we believe it is necessary to adjust the algorithm to adapt its metrics with different distributions as well. Future works may also include applying stability for evaluating edges behavior, and finding the distance between the descriptive reference process model and the individual patient pathways. This could be useful for diagnosing the causes behind the detected deviations. Furthermore, integrating hospital information systems with location data for extracting the decision points could be a valid target for future research works. This would help to extract an executable semantic which could be used for the simulation of patient pathways as business process models.

Declaration of Competing Interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

Appendix A

Algorithm 1: Stable heuristic miner.

Result: The statistically stable state: $UCL, LCL, CL, Sample.Attr$
 Input $Footprint.Matrix$; $m = \text{length}(Activity.Set)$
 $Total.Observations = \text{sum}(\text{rowSums}(Footprint.Matrix \neq 0))$
 $\bar{n} = (Total.Observations/m)$ \triangleright average size of the population
 \triangleright get the size and mean of each activity
for i in $Footprint.Matrix[1:n]$ **do**
 Sample.Attr = data.frame
 (size=rowSums(Footprint.Matrix[i,j] !=0,
 $\bar{x}_i = \text{rowMeans}(Footprint.Matrix[i,j] \neq 0,$
 $\sigma = \text{rowStandardDeviations}(Footprint.Matrix[i,j] \neq 0);$
 mutate(Sample.Attr, C_{4n}) \triangleright Eq. 5
end
 Sample.Attr[size, $\sigma, C_{4n}, \bar{x}_i$] \triangleright structure of Sample.Attr
for i in $Sample.Attr\$\bar{x}_i$ **do**
 | $\bar{x} = (\text{sum}(Sample.Attr[size * \bar{x}_i]) / \text{sum}(size))$
end
for i in $Sample.Attr[1:n]$ **do**
 | mutate(Sample.Attr, σ_i / C_{4n_i})
end
for i in $Sample.Attr[1:n]$ **do**
 | $\hat{\sigma} = ((1/m) * \text{sum}(\sigma_i / C_{4n_i}))$ \triangleright Eq. 6
end
 $A_{3\bar{n}} = 3 / C_{4\bar{n}} * \text{sqrt}(\bar{n})$
 $CL = \bar{x}$ \triangleright devise the thresholds
 $UCL = \bar{x} + (A_{3\bar{n}} * C_{4\bar{n}} * \hat{\sigma})$
 $LCL = \bar{x} - (A_{3\bar{n}} * C_{4\bar{n}} * \hat{\sigma})$

Appendix B

Algorithm 2: Stable heuristic miner.

Result: Descriptive reference process model (as a directed graph)
 Input $UCL, LCL, CL, Sample.Attr, Activity.Set$
 \triangleright considering the average of relative frequencies for each sample
for i in $Sample.Attr[\bar{x}_i]$ **do**
 Unstable.Activities = $\bar{x}_i \leq LCL$;
 Stable.Activities = $LCL < \bar{x}_i < UCL$;
 Hot.Zones = $UCL \leq \bar{x}_i$
end
 \triangleright select the nodes
 Stable.Nodes = match(Stable.Activities, Activity.Set)
 Hot.Nodes = match(Hot.Zones, Activity.Set, Color.Attr = "red")
 All.Nodes = combine(Stable.Nodes, Hot.Nodes)
 \triangleright select the edges and devise the graph
 edges = match(Footprint.Matrix, All.Nodes)
 devise.graph(All.Nodes, edges)

CRedit authorship contribution statement

Sina Namaki Araghi: Methodology, Conceptualization, Software, Validation, Writing – original draft, Visualization, Data curation. **Franck Fontanili:** Writing – review & editing, Supervision, Resources. **Elyes Lamine:** Writing – review & editing. **Uche Okongwu:** Writing – original draft, Writing – review & editing. **Frederick Benaben:** Supervision, Writing – review & editing, Project administration.

References

- Araghi, S. N., Fontanili, F., Lamine, E., Salatge, N., Lesbegueries, J., Pouyade, S. R., et al. (2018). A conceptual framework to support discovering of patients' pathways as operational process charts. In *2018 IEEE/ACS 15th international conference on computer systems and applications (AICCSA)* (pp. 1–6). IEEE.
- Augusto, A., Conforti, R., Dumas, M., Rosa, M. L., Maggi, F. M., Marrella, A., et al. (2018). Automated discovery of process models from event logs: Review and benchmark. *IEEE Transactions on Knowledge and Data Engineering*, 1. <https://doi.org/10.1109/TKDE.2018.2841877>.
- Augusto, A., Conforti, R., Dumas, M., Rosa, M. L., Maggi, F. M., Marrella, A., et al. (2019). Automated discovery of process models from event logs: Review and benchmark. *IEEE Transactions on Knowledge and Data Engineering*, 31(4), 686–705. <https://doi.org/10.1109/TKDE.2018.2841877>.
- (2009). From system complexity to emergent properties. In M. Aziz-Alaoui, & C. Bertelle (Eds.), *Understanding complex systems*. Berlin Heidelberg: Springer-Verlag. <https://doi.org/10.1007/978-3-642-02199-2>.
- Barany, I., Vu, V., et al. (2007). Central limit theorems for Gaussian polytopes. *The Annals of Probability*, 35(4), 1593–1621.
- Bezerra, F., & Wainer, J. (2013). Algorithms for anomaly detection of traces in logs of process aware information systems. *Information Systems*, 38(1), 33–44. <https://doi.org/10.1016/j.is.2012.04.004>.
- Brender, J. (2006). 7- descriptions of methods and techniques. In J. Brender (Ed.), *Handbook of evaluation methods for health informatics* (pp. 73–225). Burlington: Academic Press. <https://doi.org/10.1016/B978-012370464-1.50007-1>.
- Chapela-Campa, D., Mucientes, M., & Lama, M. (2019). Mining frequent patterns in process models. *Information Sciences*, 472, 235–257. <https://doi.org/10.1016/j.ins.2018.09.011>.
- Corchado, J. M., Bajo, J., de Paz, Y., & Tapia, D. I. (2008). Intelligent environment for monitoring Alzheimer patients, agent technology for health care. *Decision Support Systems*, 44(2), 382–396. <https://doi.org/10.1016/j.dss.2007.04.008>.
- De Cnudde, S., Claes, J., & Poels, G. (2014). Improving the quality of the heuristics miner in ProM 6.2. *Expert Systems with Applications*, 41(17), 7678–7690. <https://doi.org/10.1016/j.eswa.2014.05.055>.
- De San Pedro, J., Carmona, J., & Cortadella, J. (2015). Log-based simplification of process models. In H. R. Motahari-Nezhad, J. Recker, & M. Weidlich (Eds.), *Business process management* (pp. 457–474). Cham: Springer International Publishing.
- Dumas, M., La Rosa, M., Mendling, J., Reijers, H. A., et al. (2013). *Fundamentals of business process management: vol. 1*. Springer.
- Estrada-Torres, B., Camargo, M., Dumas, M., García-Bañuelos, L., Mahdy, I., & Yerokhin, M. (2021). Discovering business process simulation models in the presence of multitasking and availability constraints. *Data and Knowledge Engineering*, 134, 101897. <https://doi.org/10.1016/j.datak.2021.101897>.
- Fernandez-Llatas, C., Lizondo, A., Monton, E., Benedi, J.-M., & Traver, V. (2015). Process mining methodology for health process tracking using real-time indoor location systems. *Sensors*, 15(12), 29821–29840. <https://doi.org/10.3390/s151229769>.
- García, C. dos S., Meincheim, A., Faria Junior, E. R., Dallagassa, M. R., Sato, D. M. V., Carvalho, D. R., et al. (2019). Process mining techniques and applications—A systematic mapping study. *Expert Systems with Applications*, 133, 260–295. <https://doi.org/10.1016/j.eswa.2019.05.003>.
- Gorban, I. I. (2014). Phenomenon of statistical stability. *Technical Physics*, 59(3), 333–340. <https://doi.org/10.1134/S1063784214030128>.
- Gorban, I. I. (2017). The statistical stability phenomenon. *Mathematical engineering*. Springer International Publishing. <https://doi.org/10.1007/978-3-319-43585-5>.
- Janssenswillen, G. (2021). Process realism. In G. Janssenswillen (Ed.), *Unearthing the real process behind the event data: the case for increased process realism*. In *Lecture notes in business information processing* (pp. 3–19). Cham: Springer International Publishing. https://doi.org/10.1007/978-3-030-70733-0_1.
- Leemans, S., Fahland, D., & van der Aalst, W. (2014). Process and deviation exploration with inductive visual miner. In L. Limonad, & B. Weber (Eds.), *BPM demo sessions 2014 (co-located with BPM 2014, Eindhoven, The Netherlands, September 20, 2014)*, CEUR-WS.org (pp. 46–50). CEUR-WS.org. BPM Demo Sessions 2014 (BPM 2014), September 10, 2014, Eindhoven, The Netherlands, BPM 2014 ; Conference date: 10-09-2014 Through 10-09-2014
- Leno, V., Polyvyanyy, A., Dumas, M., La Rosa, M., & Maggi, F. M. (2020). Robotic process mining: Vision and challenges. *Business and Information Systems Engineering*. <https://doi.org/10.1007/s12599-020-00641-4>.
- Li, G., & van der Aalst, W. M. P. (2017). A framework for detecting deviations in complex event logs. *Intelligent Data Analysis*, 21(4), 759–779. <https://doi.org/10.3233/IDA-160044>.
- Martínez-Millana, A., Lizondo, A., Gatta, R., Vera, S., Salcedo, V. T., & Fernandez-Llatas, C. (2019). Process mining dashboard in operating rooms: Analysis of staff expectations with analytic hierarchy process. *International Journal of Environmental Research and Public Health*, 16(2), 199. <https://doi.org/10.3390/ijerph16020199>. Number: 2 Publisher: Multidisciplinary Digital Publishing Institute
- Montgomery, D. C. (2007). *Introduction to statistical quality control* (8th ed.). Industrial Engineering / Manufacturing | General & Introductory Industrial Engineering | Subjects | Wiley. <https://www.wiley.com/en-us/Introduction+to+Statistical+Quality+Control2C+8th+Edition-p-9781119399308>
- Munoz-Gama, J., Martin, N., Fernandez-Llatas, C., Johnson, O. A., Sepúlveda, M., Helm, E., et al. (2022). Process mining for healthcare: Characteristics and challenges. *Journal of Biomedical Informatics*, 127, 103994. <https://doi.org/10.1016/j.jbi.2022.103994>.

- Namaki Araghi, S. (2019). *A methodology for business process discovery and diagnosis based on indoor location data : Application to patient pathways improvement*. Ecole des Mines d'Albi-Carmaux Theses. <https://tel.archives-ouvertes.fr/tel-03037128>
- Rebuge, A., & Ferreira, D. R. (2012). Business process analysis in healthcare environments: A methodology based on process mining. *Information Systems*, 37(2), 99–116. <https://doi.org/10.1016/j.is.2011.01.003>.
- Rojas, E., Munoz-Gama, J., Sepúlveda, M., & Capurro, D. (2016). Process mining in healthcare: A literature review. *Journal of Biomedical Informatics*, 61, 224–236. <https://doi.org/10.1016/j.jbi.2016.04.007>.
- Thibaud, M., Chi, H., Zhou, W., & Piramuthu, S. (2018). Internet of things (IoT) in high-risk environment, health and safety (EHS) industries: A comprehensive review. *Decision Support Systems*, 108, 79–95. <https://doi.org/10.1016/j.dss.2018.02.005>.
- Thiede, M., Fuerstenau, D., & Barquet, A. P. B. (2018). How is process mining technology used by organizations? A systematic literature review of empirical studies. *Business Process Management Journal*. <https://doi.org/10.1108/BPMJ-06-2017-0148>.
- van der Aalst, W. M. P. (2016). *Data science in action*. Berlin, Heidelberg: Springer Berlin Heidelberg. https://doi.org/10.1007/978-3-662-49851-4_1.
- van der Aalst, W. M. P., de Medeiros, A. K. A., & Weijters, A. J. M. M. (2005). Genetic process mining. In G. Ciardo, & P. Darondeau (Eds.), *Applications and theory of petri nets 2005* (pp. 48–69). Berlin, Heidelberg: Springer Berlin Heidelberg.
- Vázquez-Barreiros, B., Mucientes, M., & Lama, M. (2015). Prodigen: Mining complete, precise and minimal structure process models with a genetic algorithm. *Information Sciences*, 294, 315–333. <https://doi.org/10.1016/j.ins.2014.09.057>. Innovative Applications of Artificial Neural Networks in Engineering
- Weijters, A. J. M. M., & van der Aalst, W. M. P. (2003). Rediscovering workflow models from event-based data using little thumb. *Integrated Computer-Aided Engineering*, 10(2), 151–162. <https://doi.org/10.3233/ICA-2003-10205>.
- Weijters, A. J. M. M., Maruster, L., & van der Aalst, W. (2004). Workflow mining: Discovering process models from event logs. *IEEE Transactions on Knowledge and Data Engineering*, 16(9), 1128–1142. <https://doi.org/10.1109/TKDE.2004.47>.
- Weijters, A. J. M. M., & Ribeiro, J. T. S. (2011). Flexible heuristics miner (FHM). In *2011 IEEE symposium on computational intelligence and data mining (CIDM)* (pp. 310–317). <https://doi.org/10.1109/CIDM.2011.5949453>.
- Zhou, W., & Piramuthu, S. (2010). Framework, strategy and evaluation of health care processes with RFID. *Decision Support Systems*, 50(1), 222–233. <https://doi.org/10.1016/j.dss.2010.08.003>.