

# Healthcare ~~Pathway Discovery and Probabilistic~~ ~~Machine Learning~~ Pathway discovery and probabilistic machine learning



The corrections made in this section will be reviewed and approved by a journal production editor.

Andreas W. Kempa-Liehr<sup>a,\*</sup> [a.kempa-liehr@auckland.ac.nz](mailto:a.kempa-liehr@auckland.ac.nz), Christina Yin-Chieh Lin<sup>a</sup>, Randall Britten<sup>b,c</sup>,  
Delwyn Armstrong<sup>d</sup>, Jonathan Wallace<sup>d</sup>, Dylan Mordaunt<sup>e,f</sup>, Michael O'Sullivan<sup>a</sup>

<sup>a</sup>Department of Engineering Science, The University of Auckland, 70 Symonds St, Auckland, New Zealand

<sup>b</sup>Auckland District Health Board, 2 Park Road, Auckland, New Zealand

<sup>c</sup>was at Orion Health, 181 Grafton Rd, Auckland, New Zealand

<sup>d</sup>Waitemata District Health Board, 124 Shakespeare Rd, Auckland, New Zealand

<sup>e</sup>University of Adelaide and Flinders University, Adelaide, Australia

<sup>f</sup>was at Waitemata District Health Board, 124 Shakespeare Rd, Auckland, New Zealand

\*Corresponding author.

---

## Abstract

**Background and purpose:** Healthcare pathways define the execution sequence of clinical activities as patients move through a treatment process, and they are critical for maintaining quality of care. The aim of this study is to combine healthcare pathway discovery with predictive models of individualized recovery times. The pathway discovery has a particular emphasis on producing pathway models that are easy to interpret for clinicians without a sufficient background in process mining. The predictive model takes the stochastic volatility of pathway performance indicators into account.

**Method:** This study utilizes the business process-mining software ProM to design a process mining pipeline for healthcare pathway discovery and enrichment using hospital records. The efficacy of combining learned healthcare pathways with probabilistic machine learning models is

demonstrated via a case study that applies the proposed process mining pipeline to discover appendicitis pathways from hospital records. Machine learning methodologies based on probabilistic programming are utilized to explore pathway features that influence patient recovery time.

**Results:** The produced appendicitis pathway models are easy for clinical interpretation and provide an unbiased overview of patient movements through the treatment process. Analysis of the discovered pathway model enables reasons for longer than usual treatment times to be explored and deviations from standard treatment pathways to be identified. A probabilistic regression model that estimates patient recovery time based on the information extracted by the process mining pipeline is developed and has the potential to be very useful for hospital scheduling purposes.

**Conclusion:** This study establishes the application of the business process modelling tool ProM for the improvement of healthcare pathway mining methods. The proposed pipeline for healthcare pathway discovery has the potential to support the development of probabilistic machine learning models to further relate healthcare pathways to performance indicators such as patient recovery time.

---

**Keywords:** Healthcare pathway; Process mining; Electronic health record; Probabilistic programming

## 1. Introduction

The success of Electronic Medical Record (EMR) systems has sparked research on healthcare knowledge management [1,2], which is a domain specific topic of the broad field of knowledge management [3,4]. Healthcare knowledge management can be described as a circular process, which consists of four steps [5]: ~~data access~~, (1) data access, (2) knowledge discovery, (3) knowledge translation & interpretation, as well as (4) knowledge description, integration & sharing. Some examples of successful healthcare knowledge management might comprise care coordination for paediatric asthma [6], the identification of patients' specific care teams [7], the enhancement of cancer care coordination [8], or the analysis of interaction patterns of trauma providers [9], to name just a few.

An important role in healthcare knowledge management are healthcare pathways, which are attributed to the operational knowledge of an organization [2]. Healthcare pathways are critical for reducing clinical variability, affecting operational excellence, and maximizing health outcomes [10]. They define the execution sequence of clinical activities as patients move through a treatment process, a department, a hospital, or a wider health organization [11]. The accurate definition of healthcare pathways and patient conformance to those pathways is an issue of increasing relevance as precision medicine enables targeted approaches and diagnostic splitting. The proliferation of pathway branches is exponential, and pathways are increasingly non-linear.

Most healthcare pathways result from clinician-led practice rather than explicit pathway design via a consensus model and systems approach. In addition, healthcare pathways *shift* dynamically as steps in the pathway are altered or resources change along the pathway. If no explicit redesign of pathways is performed, then the providers of the pathways (and its associated resources) may be unaware of the change [12]. Pathway

discovery (identification of pathways without a priori knowledge), conformance analysis (including gaps in care and clinical variability) and pathway enrichment (enrichment of a priori models with additional event data) are critical for healthcare services now, and into the future [13]. Past studies have shown that there is potential for informative healthcare pathways to be extracted from hospital health records [14–16]. Furthermore, an efficient workflow based on Business Process Modelling utilizing the process-mining software package ProM [17] has been established [18,19]. This systematic healthcare pathway mining method supports explicit design and conformance analysis of concise and comprehensible healthcare pathway models.

However, a field of research, which remains understudied despite the successes of machine learning in the context of precision driven medicine [20], is the application of healthcare pathways for predicting individualized health outcomes like recovery times. A possible explanation for this observation might be the fact, that health outcomes, which are quantified on a time domain, do exhibit in many cases fat tail distributions with dominant modes, such that classical point predictions are rarely better than predicting the respective expectation value. An alternative to regression models, which provide point predictions or prediction intervals, are generative models [21], which predict posterior predictive distributions by applying probabilistic machine learning [22].

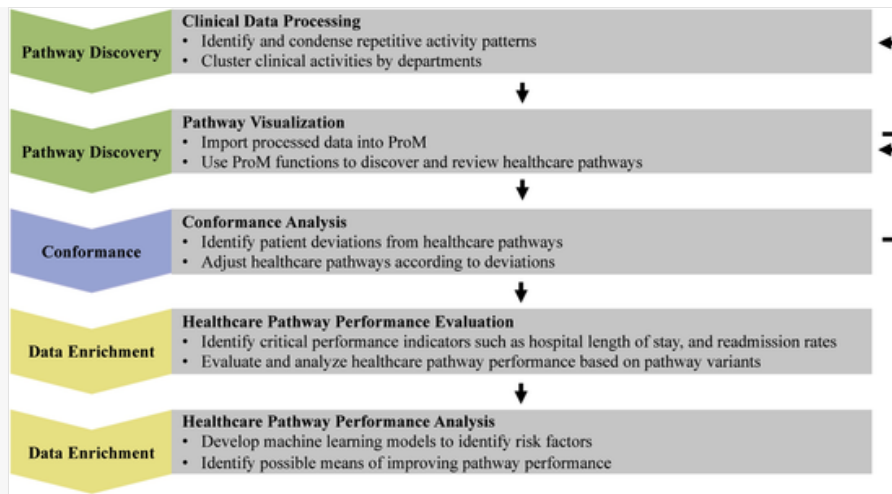
In this paper, we are presenting a case study, which uses the pathway mining software ProM for discovering the appendectomy pathways from EMRs of the North Shore Hospital in Auckland, New Zealand. The pathway discovery has a particular emphasis on producing pathway models that are easy to interpret for clinicians without a sufficient background in process mining. The main contribution of this paper is the application of probabilistic machine learning for predicting patient specific recovery times after appendectomy using pathway information from the learned pathways and other relevant EMR data. This probabilistic model is able to replicate the dominant mode as well as the fat tail of the empirical recovery time distribution and has the potential to be developed into an advanced planning tool.

The paper is organized as follows: The healthcare pathway mining methodology is described in Sec. 2. The appendectomy case study (Sec. 3) starts with a description of the EMRs used for the case study, discusses the results of the appendicitis pathway discovery, and analyses pathway specific performance indicators. The probabilistic machine learning model is introduced and evaluated in Sec. 4. The paper closes with conclusions and an outlook in Sec. 5.

## 2. Healthcare pathway mining methodology

This section outlines the process mining pipeline designed for mining and analyzing healthcare pathways using hospital patient records. This study adopts the scientific computing practices recommended by Wilson et al. [23][23,24], [24] to ensure that all results are reproducible. The process mining pipeline consists of three major phases that correspond to each of the three main objectives (i.e. pathway discovery, conformance, and enrichment). An overview of the process mining pipeline applied for this study is shown in Fig. 1. The following sections elaborate the three phases of the proposed pipeline and their respective objectives.

Fig. 1



The process mining pipeline comprises three sections, which are subdivided into five steps: **C**linical data processing, pathway visualization, conformance analysis, healthcare pathway **e**conformance**p**erformance evaluation, and predictive analytics of healthcare pathways. The first three steps are connected in two iterative cycles.

This study uses ProM (version 6.7) as the main process mining tool.<sup>1</sup> ProM is an open-source process mining software that is effective for construction of business models from input data files. ProM is chosen for this study because it has an intuitive user interface and supports many process mining plug-ins [17]. The process mining and conformance analysis plug-ins supported by ProM are well documented. All these features of ProM make it easy for the process mining steps in the pipeline (see Fig. 1) to be repeated by users with no background in process mining. The process mining techniques used in this paper are selected based on their ease for clinical interpretation and their potential to be combined with machine learning models.

## 2.1. Healthcare pathway discovery

Healthcare pathway discovery is the first phase of the process mining pipeline. It consists of two steps: clinical data processing and pathway visualization, which are conducted iteratively until a concise model is produced. The aim is to use patient healthcare records stored in hospital information systems to design a concise pathway model that is easy for clinical interpretation. Therefore, clinical input is critical to the selection of appropriate processing methods.

Healthcare pathways generally have much higher levels of complexity than standard business processes, and unprocessed clinical data contains too many clinical variations for a clean and concise pathway model to be mined [25,26]. Each pathway variant is a unique event sequence of a complete patient trace. The ProM plug-in Explore Event Log (from the Log Enhancement package) extracts pathway variants from patient traces, and the total number of pathway variants is an indicator of the level of clinical variation between patient traces. The basic format of an example pathway variant extracted by Explore Event Log is demonstrated in Fig. 2.

Fig. 2



Basic format of an example pathway variant visualized by ProM's plug-in Explore Event Log. The example process consists of three different activities (colours). The start and stop events of each activity are indicated by separate arrows, such that overlapping activities can be easily identified.

In order to reduce healthcare pathway variations to a meaningful pathway model, the pathway variants visualized by the plug-in Explore Event Log are examined closely to determine the most suitable processing methods. There are three effective methods for reducing clinical variations without filtering patient traces:

**Cluster clinical activities** that are similar in nature so that the range of activities is reduced to a manageable size, e.g., 'Abdomen CT scan' and 'Pelvis CT scan' could be clustered into a single activity under 'CT scan'.

**Merge consecutive clinical activities** that are performed consecutively into a single activity, e.g., a patient receiving the same medication five times on the same day could be regarded as a single activity.

**Condense repetitive activity patterns** that repeat but exhibit variable cycle length. These patterns indicate an activity that must be performed periodically while the patient is waiting for a different activity to begin, e.g., lab tests to monitor a patient's condition, or medication to prevent infection. These repetitive patterns could be condensed into a single, parallel activity.

Clinical input is highly recommended at this step particularly for complex or unfamiliar healthcare pathways.

## 2.2. Healthcare pathway conformance analysis

It is very difficult for clinicians to manually track individual patient traces through the treatment process and ensure that they are conforming to standard protocols. Identifying unwarranted deviations and making the required interventions early in the process has the potential to improve health outcomes and decrease cost. Conformance analysis identifies patient deviations by comparing the pathway model to clinical data. Accuracy of the discovered healthcare pathway model is validated if the majority of the patient traces conform to the model. Patient traces rarely all follow identical pathways, so the healthcare pathway model is not expected to capture all patient traces. The objective is to discover a healthcare pathway model that captures the fundamental structure of most patient traces and detect unexpected patient deviations.

ProM offers tools for conformance analysis of healthcare pathway models: Its plug-in Inductive Visual Miner compares patient traces from input clinical data to a healthcare pathway model and indicates patient traces which are deviating from the pathway model. For this purpose, the pathway model is visualized as a process tree, which is a hierarchical map comprised of decision nodes and tasks representing clinical activities [27]. Therefore, process trees enable the identification of pathway branches throughout the healthcare pathway model (c.f. Section 3.2).

If valid patient traces deviate from a healthcare pathway model, adjustments are made to the model to improve patient conformance. A typical example might be the introduction of a new form of treatment, which has not been included into the model yet. Including these findings into the model leads to an iterative approach between pathway discovery and conformance analysis (Fig. 1). Conversely, conformance analysis can identify

where invalid patient traces deviate from the model and investigate the reason for the discrepancy, e.g., clinicians following obsolete pathways or data errors.

### 2.3. Healthcare pathway data enrichment

Data enrichment of healthcare pathways is the third phase of the discussed process mining pipeline (Fig. 1). It comprises two steps: Healthcare pathway performance evaluation and healthcare pathway performance analysis.

The main objectives of evaluating healthcare pathway performance are to understand the strengths and weaknesses of the current pathway design, and to identify potential methods of improvement. Possible indicators of healthcare pathway performance include waiting times of clinical activities, hospital length of stay, recovery time, and readmission rates [28]. Most of these indicators can be calculated or estimated using standard clinical timestamps. Postoperative Length of Stay (PLS), which is measured from leaving operating theatre to discharge, can also be considered as patient recovery time. For surgical healthcare pathways, PLS is one of the critical indicators for evaluating healthcare pathway performance [29].

Analysing the performance of healthcare pathways with respect to pathway variants and other possible influencing factors like demographics or patient specific pathway characteristics, e.g. surgery duration (SD), is the final step of the process mining pipeline (Fig. 1). Due to the fact, that most pathway performance indicators do not follow normal distributions, while exhibiting significant stochastic volatility, neither classical hypotheses tests [30], nor point-predicting machine learning models are appropriate for analysing healthcare pathways. Instead, probabilistic machine learning models [22] are used for extracting interpretable models from healthcare pathways (Sec. 4.4). For this purpose, feature engineering [31] from the patients' pathway traces (e.g. SD, pathway variant), demographics (e.g. age), as well as medical documentation like written diagnosis, time series, or images become important. In order to demonstrate this approach, the following case study discusses a probabilistic machine learning model for PLS, which takes into account pathway variants (Sec. 3.2.2), as well as demographics, and SD (Sec. 4.4).

## 3. Case study: Appendicitis healthcare pathways

This section discusses the healthcare pathway mining and analysis process for an appendicitis case study. For this purpose, two years' worth of data from 2015 to 2017 on 448 appendicitis patients have been analysed. This case study is selected because clinicians confirmed it is a relatively simple surgical pathway with clear start and end points.


### 3.1. Data description

The patient records for this case study were collected from North Shore Hospital in Auckland, New Zealand. The electronic patient records in North Shore Hospital are stored using the Radiology Information System (RIS) and the patient administration system iPM. The extracted data were de-identified and an ethics approval for this research was obtained. All data sets collected from the hospital's information system on appendicitis patients are summarized in Table 1. Theatre encounter is the system ID used to identify patient traces, and clinical activities are categorized by clinical departments (e.g. radiology, pharmacy). A column labelled 'Pre-

Op/Post-Op' indicates whether a clinical activity is performed before or after surgery, and contents of this column are appended as prefixes to the clinical activities during data processing. Processed data sets are imported into ProM for pathway discovery and conformance analysis.

alt-text: Table 1

Table 1

 The table layout displayed in this section is not how it will appear in the final version. The representation below is solely purposed for providing corrections to the table. To preview the actual presentation of the table, please view the Proof.

Description of all data sets collected from North Shore Hospital for the appendicitis case study. [Instruction: The table layout needs to be corrected, because the first column is much too broad. See screenshot from Proof.]

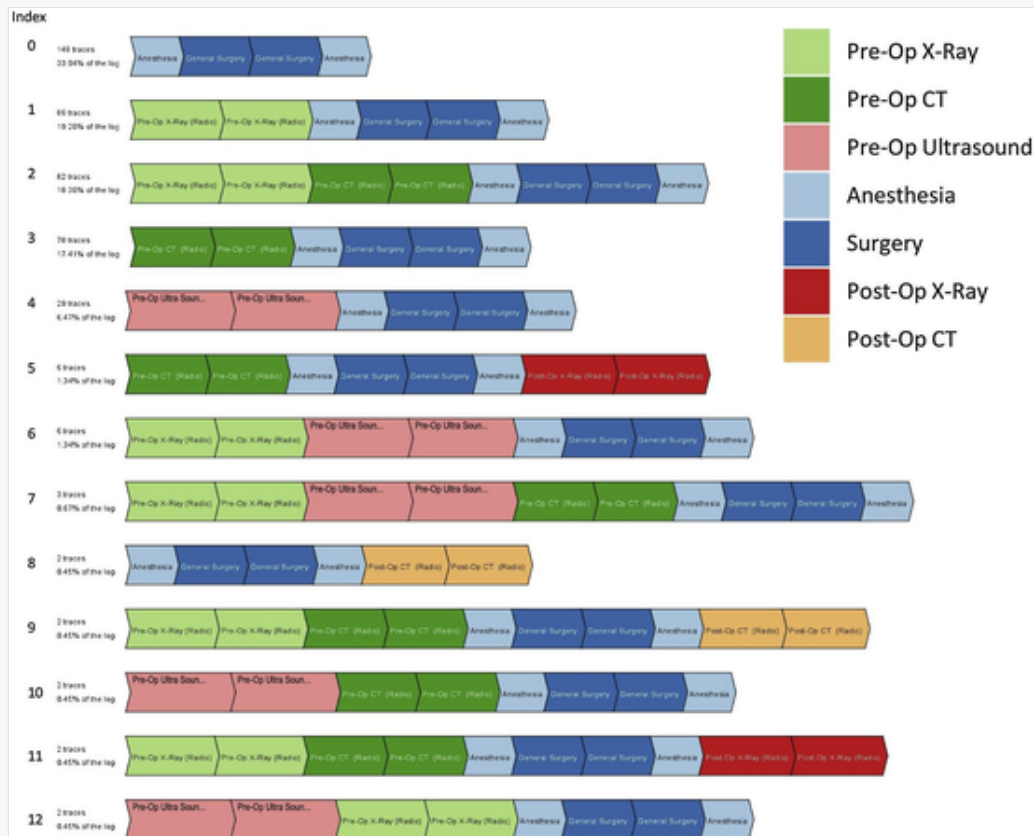
EMR	Columns of <del>Interest</del> <del>Data Type</del> <del>Acute Theatre</del> <del>Theatre Encounter</del> <del>String</del> <del>Surgery Start Time</del> <del>Datetime</del> <del>Surgery End Time</del> <del>Datetime</del> <del>Into Theatre Time</del> <del>Datetime</del> <del>Out of Theatre Time</del> <del>Interest</del>	Data type
Acute theatre	• Theatre encounter	String
	• Surgery start time	Datetime
	• Surgery end time	Datetime
	• Into theatre time	Datetime
	• <u>Out of theatre time</u>	Datetime
General <del>Surgery</del> <del>Theatre Encounter</del> <del>String</del> <del>Admission Time</del> <del>Datetime</del> <del>Discharge Time</del> <del>Surgery</del>	• Theatre encounter	String
	• Admission time	Datetime
	• <u>Discharge time</u>	Datetime
Radiology/ <del>Pharmacy</del> <del>Theatre Encounter</del> <del>String</del> <del>Anesthesia</del> <del>Clinical Activity</del> <del>String</del> <del>Clinical Activity Start Time</del> <del>Datetime</del> <del>Pre-Op/Post-Op</del> <del>pharmacy/anaesthesia</del>	• Theatre encounter	String
	• Clinical activity	String
	• Clinical activity start time	Datetime
	• <u>Pre-op/post-op</u>	String
Patient	• Theatre <del>Encounter</del> <del>String</del> <del>Patient</del> <del>Encounter</del>	String
	• <u>Patient age</u>	Integer

### 3.2. Appendicitis pathway discovery and conformance analysis



The appendicitis pathway model, which has been generated by ProM's plugin Explore Event Log, is shown in Fig. 3 in the form of a pathway variant plot. The pathway variants are extracted without activities related to medication (i.e. preoperative and postoperative cefuroxime/metronidazole) because the clinicians confirmed that antibiotics are usually taken while the patient is waiting for surgery or discharge. The duration of these activities are therefore highly variable and result in a high number of unique pathway variants.

Fig. 3



Appendicitis pathway variant plot auto-generated by ProM's plug-in Explore Event Log. The plot visualizes the sequences of start and stop events for the different pathways. For the purpose of readability the legend was added and the statistics on the left are repeated in Table 2. The top four variants account for approximately 88% of the patient traces. They are modelled as one-hot encoded features V0, V1, V2, and V3 in Section 4, while pathway variants V4–V12 together represent the base model.

The pathway variant plot visualizes the 13 pathway variants of the appendicitis model sorted from the most common pathway (index 0) to the least common pathway (index 12). The top four variants account for approximately 88% of the patient traces (Table 2). All clinical activities are represented by a start event and a stop event. The activities are colour coded such that the same colour refers to the same clinical activity. The most common pathway variant (index 0) only consists of an [a](#)esthesia and surgery, while the second most common variant (index 1) also includes preoperative X-ray. The pathway variant indices are used in Section 4 as one-hot encoded feature of the probabilistic machine learning model.



Table 2

*i* The table layout displayed in this section is not how it will appear in the final version. The representation below is solely purposed for providing corrections to the table. To preview the actual presentation of the table, please view the Proof.

Statistics of appendicitis patient traces shown in Fig.3.

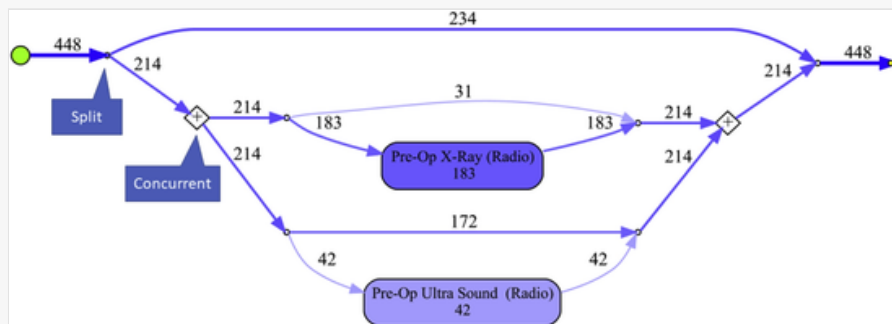
Variant	0	1	2	3	4	5	6
Patients	148	86	82	78	29	6	6
Percentage	33.04%	19.20%	18.30%	17.41%	6.47%	1.34%	1.34%

Variant	7	8	9	10	11	12
Patients	3	2	2	2	2	2
Percentage	0.67%	0.45%	0.45%	0.45%	0.45%	0.45%

The first stage of the appendicitis pathway model visualized by Inductive Visual Miner is shown in Fig. 4. Unlike the pathway variants, this appendicitis pathway model incorporates activities representing antibiotics. The model indicates that 42 patients perform ultrasound and 183 patients perform X-ray upon admission.

Fig. 4



First stage of the appendicitis pathway model generated by ProM. The following stages have been omitted for the purpose of readability. The process indicates that 234 patients do not have any preoperative imaging diagnostics, while 214 patients enter the imaging diagnostics branch. Please refer to Leeman's manual on Inductive Visual Miner for details on the model notations [32].

While the complex process tree notation of ProM's Inductive Visual Miner plug-in is optimal for detailed analysis, it has been reformulated under new notations for easy clinical interpretation. The new model notations are summarized in Table 3, and the reformulated appendicitis pathway model is shown in Fig. 5. The section of the model labelled as 'Stage 1' in Fig. 5 corresponds to the process tree shown in Fig. 4. This is the final pathway model that has been compiled based on clinical input to account for valid patient deviations, and

all patient traces conform to the updated pathway. The most common pathway variant (index 0) shown in Fig. 3 corresponds to the horizontal path from start to finish in Fig. 5. Feedback from medical experts confirmed that the simplified appendicitis pathway model (Fig. 5) correctly represents the As-Is process, while being easily interpretable.

alt-text: Table 3

Table 3

*i* The table layout displayed in this section is not how it will appear in the final version. The representation below is solely purposed for providing corrections to the table. To preview the actual presentation of the table, please view the Proof.

Definitions of new pathway model notations.





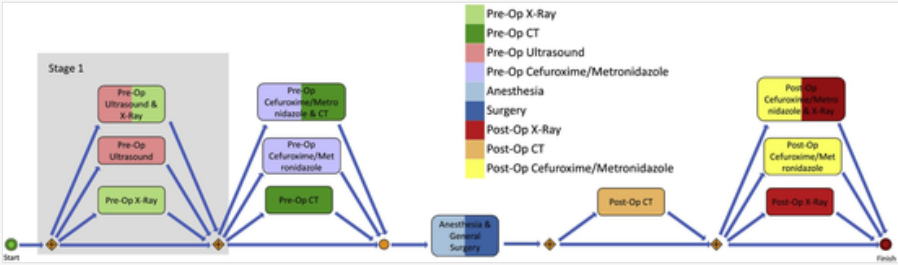
Notation	Symbol	Definition
Orange <del>Point</del> diamond		Decision <del>Point</del> point, indicating exclusive choice
<del>Orange C</del> Orange connector		Pathway <del>Connection Point</del> Green C connection point
<u>Green</u> connector		Starting <del>Point</del> Red C point
<u>Red</u> connector		Finishing <del>Point</del> point

Fig. 5

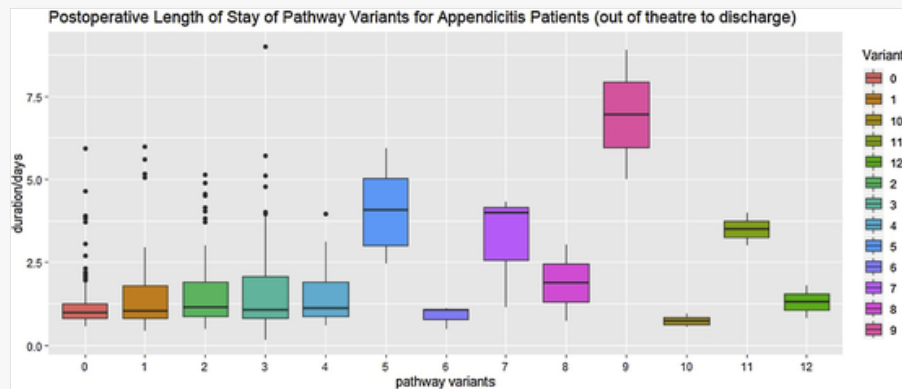


Appendicitis pathway model using nomenclature of Table 3. All preoperative and postoperative activities belong to radiology and pharmacy departments. Preoperative antibiotics are taken in the second stage of the treatment pathway.

3.3. Length of stays of appendicitis pathway variants

Length of stays in hospital are analyzed based on the 13 identified appendicitis pathway variants shown in Fig. 3. Postoperative length of stay, measured from leaving theatre to discharge, for each of the 13 appendicitis pathway variants are shown in Fig. 6. Pathway variant 9 has the longest postoperative length of stay. Pathway variants 5, 7 and 11 also have relatively long postoperative length of stays. Only pathway variant 7 does not include any postoperative activities. The number of patient traces that follow each pathway variant is limited, and more samples are required to identify rate-determining steps.

Fig. 6



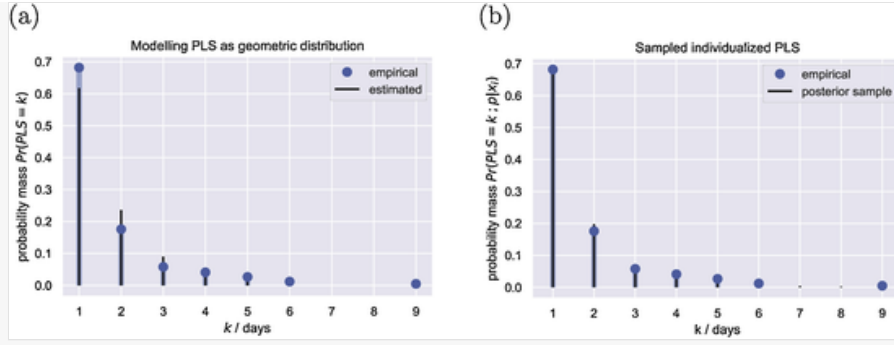
Postoperative length of stay of the 13 appendicitis pathway variants, measured from leaving theatre to discharge.

Probabilistic machine learning models are used to further investigate factors influencing postoperative length of stay (Sec. 4.4.4).

## 4. Probabilistic machine learning model for postoperative length of stay of appendicitis patients

The following section investigates the question of whether the pathway variants of the appendicitis case study are relevant features or covariates for explaining the stochastic volatility of postoperative Length of Stay (PLS, Fig. 7a), which is measured from leaving operating theatre to discharge. This is quite a challenging task, because the individual healing process is expected to depend on personal factors like age [33] as well as the individual severity of the appendicitis inflammation, which in general is unknown at this stage of the data analysis, but might be captured by proxies like surgery duration (SD).

Fig. 7



Modelling PLS as geometric distribution without patient specific model (a) and with individualized posterior (b) from model B (Eq. (4)). Sampling from patient specific posterior distributions replicates the observed PLS quite well.

For the purpose of this analysis, patients' traces with complete demographic information (415 samples) are selected and the patients' *PLS*, which is calculated as the time interval between leaving operating theatre and discharge, is converted into the number of days after surgery. Therefore, a *PLS* of one corresponds to a patient who has been discharged the day after surgery. Three of the  $N = 415$  patient traces had a *PLS* of zero, because the surgery ended shortly after midnight and the patients were discharged on the same day. In order to simplify our analysis, the interpretation of *PLS* was broadened to the number of night rests after surgery, such that the *PLS* of these three patients could be projected to one.

Analysing the *PLS* of the 415 appendicitis patients reveals a histogram which closely resembles a Geometric distribution (Fig. 7a)

$$\Pr(PLS = k) = \text{Geom}(k; p) = (1 - p)^{k-1} p^k,$$

(1)

parameterized by probability  $p$  of being discharged on the  $k$ th day, respectively  $k$  nights of rest after surgery. However, estimating  $p$  as inverse mean of the observed *PLS* to [Instruction: Please make sure that Eq (2) is correctly rendered in the PDF, which is not the case right now.]

#### Previous Version

$$\hat{p} = \frac{N}{\sum_{i=1}^N PLS_i} \approx 61.8\%,$$

#### Updated Version

$$\hat{p} = \frac{N}{\sum_{i=1}^N PLS_i} \approx 61.8\%,$$

(2)

reveals that the average discharge probability  $\hat{p}$  does not generalize well over the cohort (Fig. 7a), because the number of patients being discharged on the first day after surgery is underestimated and the number of patients being discharged on the second and third day after surgery is overestimated.

For the following probabilistic machine learning models, the individualized discharge probability  $p_i$  is modelled as a generalized linear model (GLM) using the inverse logistic function (logit) as a link function and the geometric distribution for generating the likelihood. For the following analysis, we are comparing two different models of discharge probability  $p_i$  as

$$\text{logit}(p_i) \sim SD_i + \log(\text{age}_i) + V_{0,i} + V_{1,i} + V_{2,i} + V_{3,i} + V_{0,i} \times \log(\text{age}_i) + V_{1,i} \times \log(\text{age}_i) + V_{2,i} \times \log(\text{age}_i) + V_{3,i} \times \log(\text{age}_i), \quad (3)$$

which we refer to as *model A*, and

$$\text{logit}(p_i) \sim SD_i + \log(\text{age}_i) + V_{0,i} + V_{1,i} + V_{2,i} + V_{3,i} + V_{0,i} \times \log(\text{age}_i) + V_{1,i} \times \log(\text{age}_i), \quad (4)$$

which we refer to as *model B*. These two models have been selected because they were the most credible based on the widely applicable information criterion (WAIC) and leave-one-out cross-validation (LOO) statistics [34]. The explanatory variables  $V_{j,i}$  are categorical with

$$V_{j,i} = \begin{cases} 1 & : V_{j,i} = j, \\ 0 & : \text{otherwise.} \end{cases} \quad (5)$$

Therefore, the case  $V_{0,i} = V_{1,i} = V_{2,i} = V_{3,i} = 0$  corresponds to variants V4–V12 of Fig. 3. Models A and B make the simplifying assumption that the individualized probability of discharge  $Pr(PLS_i = k | x_i)$  can be estimated from

$$x_i = (SD_i, \log(\text{age}_i), V_{0,i}, V_{1,i}, V_{2,i}, V_{3,i}). \quad (6)$$

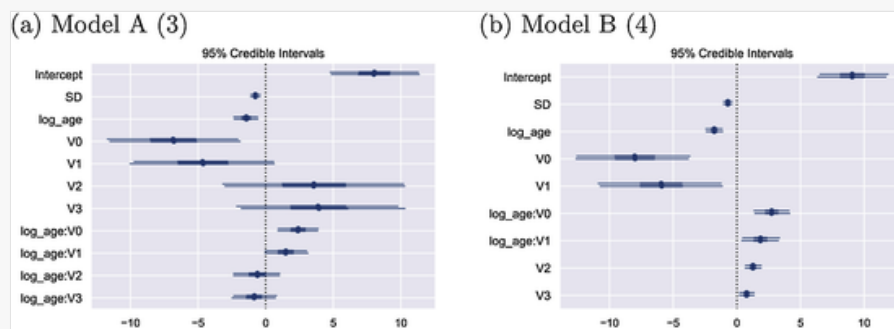
being available at the end of surgery, when future treatment steps and therefore the complete pathway variant for the respective  $i$ th patient have been decided on: [Instruction: Please make sure that both lines of Eq. (7) are aligned at their equal signs.]

$$\begin{aligned} Pr(PLS_i = k | x_i) &= \text{Geom}(k; p(x_i)) \\ &= (1 - p(x_i))^{k-1} (p(x_i))^k \end{aligned} \quad (7)$$

with  $p(x_i)$  being modelled either by Eq. (3) or (4).

The statistical model has been fitted using PyMC3 [35] version 0.24.2 using the No-U-Turn Sampler (NUTS) with 20,000 samples, 2000 tuning steps, 2 chains, and an acceptance rate of 90%. The 95% credible intervals of the estimated model coefficients are shown in Fig. 8. The Gelman–Rubin convergence statistic  $R_{hat}$  is close to one for all coefficients (Table 4). However, for Model A the coefficients of the interaction terms  $\log(age_i) \times V_2$  and  $\log(age_i) \times V_3$  are close to zero with negative expectation values but a 25% and 15% probability of being larger than zero (Table 4). Therefore, we have decided to focus the following analysis of the fitted model on Eq. (4) for which all credible intervals exclude zero (Fig. 8b). However, this does not change the overall trend depicted in Fig. 9. It is important to note that both the coefficients of  $SD$  and  $\log(age)$  are negative. This means that for the base model of variants V4–V12 as well as variants V2 and V3, the probability of discharge decreases with increasing  $SD$  and age. Compared to the base model, Variants V2 and V3 have a higher discharge probability because both coefficients are positive. The situation is more complicated for variants V0 and V1, because their coefficients are negative but the coefficients of their interaction terms with  $\log(age)$  are positive, which basically counterbalances the effect of  $\log(age)$ .

Fig. 8



Credible intervals for models A (3) and B (4).

alt-text: Table 4

Table 4

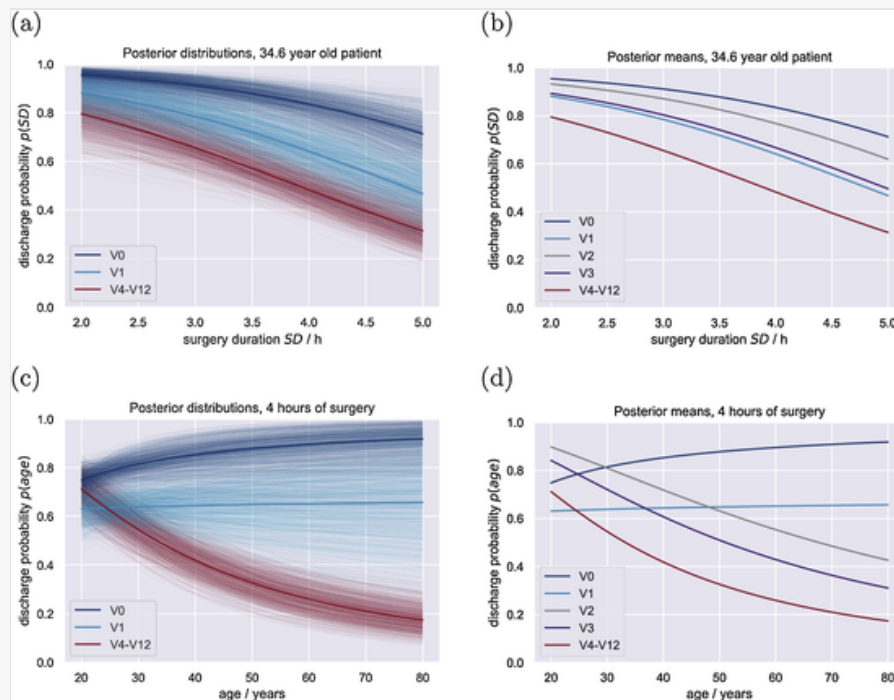
*i* The table layout displayed in this section is not how it will appear in the final version. The representation below is solely purposed for providing corrections to the table. To preview the actual presentation of the table, please view the Proof.

Coefficients of model B.  $R_{hat}$  is the potential scale reduction factor. A value of one indicates convergence.

	mMean	hpd 2.5	hpd 97.5	Rhat
Intercept	9.09	6.48	11.85	0.999988

SD	<del>-0.71</del> <del>1.059</del> -0.71	-1.059	<del>0.371</del>	0.999987
log age	<del>-1.781</del> <del>2.46</del> <del>1.154</del> <del>0.999988</del> <del>V0-8.04</del> <del>12.56</del> <del>3.700</del> <del>0.999980</del> <del>V1-5.95</del> <del>10.82</del> -1.781	-2.46	-1.154	0.999988
V0	-8.04	-12.56	-3.70	0.999980
V1	-5.95	-10.82	<del>1.186</del>	1.000013
log age: V0	2.75	1.365	4.12	0.999980
log age: V1	1.864	0.423	3.32	1.000013
V2	1.267	0.646	1.901	1.000062
V3	0.766	0.199	1.363	1.000046

Fig. 9



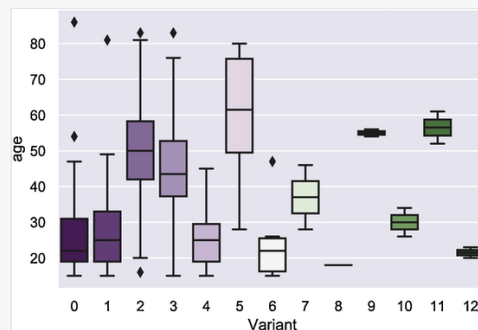
Posterior distributions and corresponding means for different pathway variants of model B (4). Trends with respect to SD and age are similar for model A (4).

This is visualized in Fig. 9, which shows the posterior distributions of discharge probability  $p_i$  as functions of SD (Fig. 9a) and age (Fig. 9c). While the discharge probability decreases with increasing SD for all pathway variants (Fig. 9b), the effect is different for age. The discharge probability  $p_i$  decreases for Variants V2, V3, and V4-V12 with increasing age of the patients, while its expectation value is constant for V1 and even increases by nearly 10% for V0. Note, that the corresponding graphics for model A (3), which are not shown here, resemble nearly the same correlations. From the clinical perspective, the difference in effects of age with respect to the pathway variants might be related to the fact that the most common pathway variants V0 and V1 are predominantly followed by younger patients (Fig. 10) and possibly older patients, for whom no



complications are expected. However, the uncertainty of discharge probabilities for V1 is significantly larger compared to V0.

Fig. 10



Boxplot of patients' age for pathway variants of the appendicitis model.

## 5. Conclusion and outlook

Healthcare pathways are critical for maintaining quality of care and improving health outcome for all patients, but there is no consensus on a healthcare pathway mining pipeline suitable for hospital implementation that supports pathway discovery from hospital health records. Business process modelling methods are used to design a process mining pipeline that produces concise and comprehensible healthcare pathway models from hospital records, and supports conformance analysis and enrichment of the discovered pathways. The proposed process mining pipeline successfully constructs concise pathway models for the appendicitis case study. The produced healthcare pathway models are easy for clinical interpretation and provide an unbiased overview of real patient traces through the treatment process. Probabilistic machine learning models for predicting postoperative length of stay, using information extracted by the process mining pipeline, **isare** showing promising results. This means that the proposed mining pipeline has the potential to support the development of machine learning models to further relate healthcare pathways to performance indicators. This study has established the use of business process modelling methods for the improvement of healthcare pathway mining methods, and there is value in investigating the capabilities of other business process modelling tools for healthcare pathway mining purposes.

[Instruction: This box is incorrectly rendered in the PDF version. The words "Summary points" are overlapping the words "Healthcare pathways".]Summary points

What was already known:

- Healthcare pathways are critical for reducing clinical variability, affecting operational excellence, and thereby maximizing health outcomes.
- Most healthcare pathways result from clinician-led practice rather than explicit pathway design via a consensus model and systems approach.

- ProM has been established as **an** efficient tool for healthcare **pathway** mining.

What this study adds:

- Pathway variants and other features from event sequences can be used as covariates for probabilistic machine learning models in order to explain the observed stochastic volatility of pathway performances.
- Geometric regression can generate posterior predictive distributions of individualized recovery times.
- Probabilistic machine learning models enable the prediction of patient specific discharge probabilities, which lead to patient specific postoperative length of stay distributions.

## Author contributions

CL performed data analysis and processing using ProM, and created the new pathway visualization method. AKL developed the probabilistic machine learning component of the work. MS and RB contributed during the conceptualization of the project. MS, RB and AKL contributed high level steering and guidance during the pathways discovery phase of the work. DA assisted with data acquisition and understanding the context of the data. JW provided guidance on clinical subject matter. DM and MS scoped the project, DM secured the funding, and MS provided overall leadership of the project.


## Declaration of competing interests

The authors declare that they have no competing interests for this study.

## Acknowledgements

The authors like to thank Patrick Gladding and Ilze Ziedins for fruitful discussions. This research has been funded by the Precision Driven Health Research Partnership (project number 1209).

## References

 The corrections made in this section will be reviewed and approved by a journal production editor. The newly added/removed references and its citations will be reordered and rearranged by the production team.

[1] Bali R.K., Dwivedi A.N. (Eds.), Healthcare Knowledge Management. Issues, Advances, and Successes, Health Informatics Series, Springer, New York, 2007.

[2] Haas P., Wissensmanagement in der Medizin, Forum 31 (1) (2016) 28–32, doi:10.1007/s12312-015-0012-6.

[3]

Nonaka I., Takeuchi H., ~~The knowledge-creating company~~[The Knowledge-creating Company](#), Oxford University Press, New York, 1995.

- [4] Beier H., Schmidt U., Klett D. (Eds.), Wissensmanagement Beflügelt, Akademische Verlagsgesellschaft, Heidelberg, 2015.
- [5] Frize M., Walker R.C., Catley C., ~~Healthcare knowledge management: Knowledge management in the perinatal care environment~~[Healthcare knowledge management: knowledge management in the perinatal care environment](#), in: Bali R.K., Dwivedi A.N. (Eds.), Healthcare Knowledge Management. Issues, Advances, and Successes, Health Informatics Series, Springer, New York, 2007, pp. 232–259.
- [6] Janevic M.R., Baptist A.P., Bryant-Stephens T., Lara M., Persky V., Ramos-Valencia G., Uyeda K., Hazan R., Garrity A., Malveaux F.J., Effects of pediatric asthma care coordination in underserved communities on parent perceptions of care and asthma-management confidence, ~~Journal of Asthma. Asth.~~[J. Asth.](#) 54 (5) (2017) 514–519, doi:10.1080/02770903.2016.1242136.
- [7] He S., Gurr G., Rea S., Thornton S.N., Characterizing the structure of a patient's care team through electronic encounter data analysis, in: Indra Neil Sarkar A.G., Mazzoncini de Azevedo Marques P. (Eds.), MEDINFO 2015: eHealth-Enabled Health: Proceedings of the 15th World Congress on Health and Biomedical Informatics, Vol. 216 of Studies in Health Technology and Informatics, IOS Press, Amsterdam, 2015, pp. 21–24, doi:10.3233/978-1-61499-564-7-21.
- [8] Mayer D.K., Deal A.M., Crane J.M., Chen R.C., Asher G.N., Hanson L.C., Wheeler S.B., Gerstel A., Green M.A., Birken S.A., Rosenstein D.L., ~~Using survivorship care plans to enhance communication and cancer care coordination: Results of a pilot study~~[Using survivorship care plans to enhance communication and cancer care coordination: results of a pilot study](#), ~~Onco Nurs Forum~~[Oncol. Nurs. Forum](#) 43 (5) (2016) 636–645, doi:10.1188/16.ONF.636-645.
- [9] Chen Y., Patel M.B., McNaughton C.D., Malin B.A., Interaction patterns of trauma providers are associated with length of stay, ~~Journal of the American Medical Informatics Association. Am. Med. Inform. Assoc. J. Am. Med. Inform. Assoc.~~ 25 (7) (2018) 790–799, doi:10.1093/jamia/ocy009.
- [10] Lin F.-R., Chou S.-C., Pan S.-M., Chen ~~Y.-M.~~, Mining time dependency patterns in clinical pathways, ~~International Journal of Medical Informatics. J. Med. Inform.~~[Int. J. Med. Inform.](#) 62 (1) (2001) 11–25, doi:10.1016/S1386-5056(01)00126-5.
- [11] Huang Z., Dong W., Ji L., He C., Duan H., Incorporating comorbidities into latent treatment pattern mining for clinical pathways, ~~Journal of Biomedical Informatics. Biomed. Inform. J. Biomed. Inform.~~ 59 (2016) 227–239, doi:10.1016/J.JBI.2015.12.012.
- [12] Zhang Y., Padman R., Wasserman L., Patel N., Teredesai P., Xie Q., ~~On Clinical Pathway Discovery from Electronic Health Record Data~~[On clinical pathway discovery from electronic](#)

- [health record data](#), ~~IEEE Intelligent Systems. Syst.~~[IEEE Intell. Syst.](#) 30 (1) (2015) 70–75, doi:10.1109/MIS.2015.14.
- [13] Baker K., Dunwoodie E., Jones R.G., Newsham A., Johnson O., Price C.P., Wolstenholme J., Leal J., McGinley P., Twelves C., Hall G., Process mining routinely collected electronic health records to define real-life clinical pathways during chemotherapy, ~~International Journal of Medical Informatics. J. Med. Inform.~~[Int. J. Med. Inform.](#) 103 (2017) 32–41, doi:10.1016/j.ijmedinf.2017.03.011.
- [14] Iwata H., Hirano S., Tsumoto S., Mining clinical pathway based on clustering and feature selection, ~~In: Brain and Health Informatics. BHI 2013, Vol. 8211 of Lecture Notes in Computer Science~~[Brain and Health Informatics. BHI 2013, Vol. 8211 of Lecture Notes in Computer Science](#), Springer, Cham, 2013, pp. 237–245, doi:~~0.1007/978-3-319-02753-1\_24~~[10.1007/978-3-319-02753-1\\_24](#).
- [15] Xu X., Jin T., Wei Z., Wang J., ~~Incorporating Topic Assignment Constraint and Topic Correlation Limitation into Clinical Goal Discovering for Clinical Pathway Mining~~[Journal of Healthcare Engineering](#)~~topic assignment constraint and topic correlation limitation into clinical goal discovering for clinical pathway mining~~[Incorporating topic assignment constraint and topic correlation limitation into clinical goal discovering for clinical pathway mining](#), J. Healthc. Eng. (2017) 1–13, doi:10.1155/2017/5208072.
- [16] Yan C., Chen Y., Li B., Liebovitz D., Malin B., Learning clinical workflows to identify subgroups of heart failure patients, ~~AMIA Annu Symp Proceedings. Symp. Proc.~~[AMIA Annu. Symp. Proc.](#) 2016 (2017) 1248–1257.
- [17] Van Dongen B.F., De Medeiros A.K.A., Verbeek H.M.W., Weijters A.J.M.M., Van Der Aalst W.M.P., ~~The ProM framework: A new era in process mining tool support~~[The ProM framework: a new era in process mining tool support](#), in: Ciardo G., Darondeau P. (Eds.), International Conference on Application and Theory of Petri Nets, Vol. 3536 of Lecture Notes in Computer Science, Springer, 2005, pp. 444–454, doi:10.1007/11494744\_25.
- [18] Mans R.S., van der Aalst W.M.P., Vanwersch R.J.B., Process Mining in Healthcare. Evaluating and Exploiting Operational Healthcare Processes, Springer Briefs in Business Process Management, Springer, Cham, 2015, doi:10.1007/978-3-319-16071-9.
- [19] Quintano Neira R.A., Hompes B.F.A., de Vries G.-J., Mazza B.F., de Almeida S.L.S., Stretton E., Buijs J.C., Hamacher S., Analysis and optimization of a sepsis clinical pathway using process mining, ~~In: International Workshop on Process-Oriented Data Science for Healthcare, International Conference on Business Process Management (BPM 2019), POS4H~~[International Workshop on Process-Oriented Data Science for Healthcare, International Conference on Business Process Management \(BPM 2019\), POS4H](#), Vienna, 2019, pp. 1–12.

[20]

Adam T., Aliferis C.F. (Eds.), Personalized and Precision Medicine Informatics: A Workflow-Based View, Springer, Cham, 2020, doi:10.1007/978-3-030-18626-5.

- [21] Bishop C.M., ~~Pattern Recognition and Machine Learning, Information science and statistics~~Pattern Recognition and Machine Learning, Information Science and Statistics, Springer, New York, 2006.
- [22] Ghahramani Z., Probabilistic machine learning and artificial intelligence, *Nature* 521 (7553) (2015) 452, doi:10.1038/nature14541.
- [23] Wilson G., Aruliah D.A., Brown C.T., Chue Hong N.P., Davis M., Guy R.T., Haddock S.H.D., Huff K.D., Mitchell I.M., Plumbley M.D., Waugh B., White E.P., Wilson P., ~~Best Practices for Scientific Computing~~Best practices for scientific computing, PLoS BiologyPLoS Biol. 12 (1) (2014) e1001745, doi:10.1371/journal.pbio.1001745.
- [24] Wilson G., Bryan J., Cranston K., Kitzes J., Nederbragt L., Teal T.K., Good enough practices in scientific computing, ~~PLOS Computational Biology: Biol.~~PLOS Comput. Biol. 13 (6) (2017) e1005510, doi:10.1371/journal.pcbi.1005510.
- [25] Huang Z., Lu X., Duan H., Fan W., Summarizing clinical pathways from event logs, ~~Journal of Biomedical Informatics. Biomed. Inform.~~J. Biomed. Inform. 46 (1) (2013) 111–127, doi:10.1016/J.JBI.2012.10.001.
- [26] Veiga G.M., Ferreira D.R., ~~Understanding Spaghetti Models with Sequence Clustering for ProM, in: spaghetti models with sequence clustering for ProM~~Understanding spaghetti models with sequence clustering for ProM, International Conference on Business Process ManagementInternational Conference on Business Process Management, Springer, Berlin, Heidelberg, 2010, pp. 92–103, doi:~~0.1007/978-3-642-12186-9\_10~~10.1007/978-3-642-12186-9\_10.
- [27] Leemans S., Fahland D., Aalst W. der, Process and deviation exploration with inductive visual miner, in: Limonad L., Weber B. (Eds.), ~~BPM Demo Sessions 2014 (Co-Located with BPM 2014, Eindhoven, The Netherlands, September 20, 2014)~~2014CEUR Workshop Proceedings, CEUR-WS.orgBPM Demo Sessions 2014 (Co-Located with BPM 2014, Eindhoven, The Netherlands, September 20, 2014), CEUR Workshop Proceedings, CEUR-WS.org, 2014, pp. 46–50.
- [28] Rotter T., Kugler J., Koch R., Gothe H., Twork S., van Oostrum J.M., Steyerberg E.W., A systematic review and meta-analysis of the effects of clinical pathways on length of stay, hospital costs and patient outcomes, ~~BMC Health Services Research 2008~~BMC Health Serv. Res. 2008 8 (265) (2008) 1–15, doi:10.1186/1472-6963-8-265.
- [29] ~~S. D. Pearson, S. F. Kleefield, J. R. The American Journal of Medicine~~S.D. Pearson, S.F. Kleefield, J.R. Soukop, E. Cook, T.H. Lee, Critical pathways intervention to reduce length of

[hospital stay, Am. J. Med. doi:10.1016/S0002-9343\(00\)00705-1.](#)

- [30] ~~S. Goodman, A dirty dozen: Twelve p-value misconceptions, Seminars in Hematology 45 (3). doi:10.1053/j.seminhematol.2008.04.003.~~ Goodman S., A dirty dozen: twelve p-value misconceptions, Semin. Hematol. 45 (3) (2008) 135–140, doi:10.1053/j.seminhematol.2008.04.003
- [31] Dong G., Liu H. (Eds.), Feature Engineering for Machine Learning and Data Analytics, Taylor & Francis, Boca Raton, FL, 2018.
- [32] Leemans S.J., ~~Inductive visual miner2017Visual Miner~~Inductive Visual Miner, 2017 <http://www.leemans.ch/publications/ivm.pdf>.
- [33] Polanczyk C., Marcantonio E., Goldman L., Rohde L., Orav J., Mangione C.M., Lee T., Impact of age on perioperative complications and length of stay in patients undergoing noncardiac surgery, ~~Annals of Internal Medicine. Intern. Med.~~Ann. Intern. Med. 134 (2001) 637–643, doi:10.7326/0003-4819-134-8-200104170-00008.
- [34] Vehtari A., Gelman A., Gabry, Practical bayesian model evaluation using leave-one-out cross-validation and waic, ~~Statistics and Computing. Comput.~~Stat. Comput. 27 (5) (2017) 1413–1432, doi:10.1007/s11222-016-9696-4.
- [35] Salvatier J., Wiecki T.V., Fonnesbeck C., Probabilistic programming in Python using PyMC3, ~~PeerJ Computer Science. Sci.~~PeerJ Comput. Sci. 2 (e55) (2016) 1–24, doi:10.7717/peerj-cs.55.

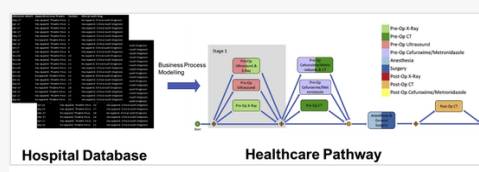
## Footnotes

### Text Footnotes

[<sup>1</sup>] <http://www.promtools.org>.

## Graphical abstract

Figure Replacement Requested



**Replacement Image:** graphical\_abstract.pdf

**Replacement Instruction:** Replace image requested

---

## Highlights

- Pathway variants and other features from event sequences can be used as co-variates for probabilistic machine learning models in order to explain the observed stochastic volatility of pathway performances.
  - Geometric regression can generate posterior predictive distributions of individualized recovery times.
  - Probabilistic machine learning models enable the prediction of patient specific discharge probabilities, which lead to patient specific postoperative length of stay distributions.
- 

## Queries and Answers

**Query:** The author names have been tagged as given names and surnames (surnames are highlighted in teal color). Please confirm if they have been identified correctly.

**Answer:** Yes

**Query:** “Your article is registered as a regular item and is being processed for inclusion in a regular issue of the journal. If this is NOT correct and your article belongs to a Special Issue/Collection please contact [r.southworth@elsevier.com](mailto:r.southworth@elsevier.com) immediately prior to returning your corrections.”.

**Answer:** The article is a regular item.

**Query:** One or more sponsor names and the sponsor country identifier may have been edited to a standard format that enables better searching and identification of your article. Please check and correct if necessary.

**Answer:** Yes

**Query:** Please provide year of publication for refs. [29,30].

**Answer:** S. D. Pearson, S. F. Kleeefield, J. R. Soukop, E. Cook, T. H. Lee, Critical pathways intervention to reduce length of hospital stay, The American Journal of Medicine 110 (3) (2001) 175–180. doi:10.1016/S0002-9343(00)00705-1.

**Query:** Highlights should only consist of 85 characters per bullet point, including spaces. The highlights provided are too long; please edit them to meet the requirement.



**Answer:** These scientific statements cannot be squeezed into 85 characters without distorting the meaning.