

Pariket: Mining Business Process Logs for Root Cause Analysis of Anomalous Incidents

Nisha Gupta

Computer Science

Indraprastha Institute of Information Technology, Delhi (IIIT-D), India

A Thesis Report submitted in partial fulfilment for the degree of

MTech Computer Science

13 February 2015

1.: Prof. Ashish Sureka (Thesis Adviser)

2. Dr. Shubhashis Sengupta (External Examiner)

2. Prof. Chetan Arora (Internal Examiner)

Day of the defense: 13 February 2015

Signature from Post-Graduate Committee (PGC) Chair:

Abstract

Process mining consists of extracting knowledge and actionable information from event-logs recorded by Process Aware Information Systems (PAIS). PAIS are vulnerable to system failures, malfunctions, fraudulent and undesirable executions resulting in anomalous trails and traces. The flexibility in PAIS resulting in large number of trace variants and the large volume of event-logs makes it challenging to identify anomalous executions and determining their root causes. We propose a framework and a multi-step process to identify root causes of anomalous traces in business process logs. We first transform the event-log into a sequential dataset and apply Window-based and Markovian techniques to identify anomalies. We then integrate the basic eventlog data consisting of the Case ID, time-stamp and activity with the contextual data and prepare a dataset consisting of two classes (anomalous and normal). We apply Machine Learning techniques such as decision tree classifiers to extract rules (explaining the root causes) describing anomalous transactions. We use advanced visualization techniques such as parallel plots to present the data in a format making it easy for a process analyst to identify the characteristics of anomalous executions. We conduct a triangulation study to gather multiple evidences to validate the effectiveness and accuracy of our approach.

I dedicate my MTech Thesis to my father-in-law Surinder Gupta who is not here to live this day. He encouraged and supported me to pursue MTech after marriage. I miss him.

Acknowledgements

Foremost, I would like to thank my advisor Prof. Ashish Sureka for the continuous support of my thesis, for his patience, motivation, enthusiasm, and immense knowledge. His guidance, invaluable constructive criticism and friendly advice helped me in all the time of research and writing of this thesis. This thesis has only been possible because of you. I could not have imagined having a better advisor and mentor for my MTech thesis.

Next, I would like to thank my fellow mate Kritika Anand for her constant support, suggestions, encouragement and help during the course of my thesis. I would also like to thank all of my friends who supported me to strive towards my goal.

Finally, a special thanks to my family. Words cannot express how grateful I am to my in-laws, my grandparents, my parents and brother for their selfless love, support and encouragement. At the end I would like to thank my husband Sagar Gupta who spent sleepless nights with me and was always my support.

Declaration

This is to certify that the MTech Thesis Report titled **Pariket: Mining Business Process Logs for Root Cause Analysis of Anomalous Incidents** submitted by **Nisha Gupta** for the partial fulfillment of the requirements for the degree of *MTech in Computer Science* is a record of the bonafide work carried out by her under my guidance and supervision at Indraprastha Institute of Information Technology, Delhi. This work has not been submitted anywhere else for the reward of any other degree.

Professor Ashish Sureka

Indraprastha Institute of Information Technology, New Delhi

Contents

List of Figures	v
List of Tables	vii
1 Research Motivation and Aim	1
2 Related Work and Novel Research Contributions	3
2.1 Related Work	3
2.2 Novel Research Contributions	5
3 Research Framework and Solution Approach	6
4 Experimental Dataset	9
5 Experimental Results	12
5.1 Sequential Dataset Conversion	12
5.2 Anomaly Detection	14
5.2.1 Window Based Technique	15
5.2.2 Markovian Based Technique	17
5.3 Data Pre-processing	19
5.4 Classification	21
5.5 Visualization	26
5.6 Triangulation Study	29
6 Limitations and Future Work	31
7 Conclusion	32
References	33

List of Figures

3.1	Architecture diagram and data processing pipeline for Pariket (Mining Business Process Logs for Root Cause Analysis of Anomalous Incidents)	7
4.1	Information Technology Infrastructure Library (ITIL) process implemented in Rabobank Group	9
4.2	Output of CA ERwin Data Modeler	10
5.1	Pareto chart showing the distribution of activities and their cumulative count. Y-axis is in logarithmic scale	12
5.2	Activities ordered according to increasing order of datetime stamp for incident id IM0000004	13
5.3	Incident id's with the sequence of activities separated by semicolon . .	14
5.4	Graphical user interface of Weka	22
5.5	Flow of classification using Weka	22
5.6	Decision tree classifier algorithm names and confusion matrices on records from Markovian and Window based technique	23
5.7	Fragments of decision tree using Weka based on anomalous incidents received from Markovian based technique	24
5.8	Fragment of decision tree using Weka based on anomalous incidents received from Window based technique	24
5.9	Decision tree using Weka based on anomalous incidents received from Markovian based technique	25

LIST OF FIGURES

5.10	Parallel coordinate plot depicting the behaviour of incident CI type affected, interaction open time, incident priority, incident open time and incident category for anomalous and non-anomalous incidents from Markovian technique	26
5.11	Parallel coordinate plot depicting the behaviour of incident CI type affected, interaction open time, incident priority, incident open time and incident category for anomalous incidents from Markovian technique . .	27
5.12	Parallel coordinate plot depicting the behaviour of incident CI type affected, interaction open time, incident priority, incident open time and incident category for anomalous incidents from Window based technique	27
5.13	Parallel coordinate plot depicting the behaviour of incident CI type affected, incident CI subtype affected, incident category, interaction CI type affected and interaction CI subtype affected for anomalous and non-anomalous incidents from Markovian technique	28
5.14	Parallel coordinate plot depicting the behaviour of incident CI type affected, incident CI subtype affected, incident category, interaction CI type affected and interaction CI subtype affected for anomalous incidents from Markovian technique	28
5.15	Parallel coordinate plot depicting the behaviour of incident CI type affected, incident CI subtype affected, incident category, interaction CI type affected and interaction CI subtype affected for anomalous incidents from Window based technique	29

List of Tables

5.1	Name, Type and Description of Some of Attributes in Merged table of Interaction_detail and Incident_detail	20
-----	---	----

Research Motivation and Aim

A Process-Aware Information System (PAIS) is a software system that manages and executes operational processes involving people, applications, and/or information sources on the basis of process models [1]. Example PAIS's are workflow management systems, case-handling systems, enterprise information systems, etc. PAIS log events and activities during the execution of a process. Process Mining is a relatively young and emerging discipline consisting of analyzing the event logs from such systems for extracting knowledge such as the discovery of runtime process model (discovery), checking and verification of the design time process model with the runtime process model (conformance analysis) and improving the business process (recommendation and extension) [2] [3]. Process mining uses data mining techniques in the context of business process management and enables the application of innovative approaches for improving the management of business processes. Process mining techniques attempt to extract non-trivial and useful information from event logs. The knowledge obtained this way can increase understanding of the processes within an organization. An event log is a collection of events. A process consists of cases or incidents. A case is a record of events that relate to a single executed process instance. Events within a case are ordered and have attributes such as activity, timestamp, actor and several additional information such as the cost. The incidents and activities in event logs can be modeled as sequential and time-series data. Anomaly detection in business process logs is an area that has attracted several researcher's attention [4] [5]. Anomalies are patterns in data that do not conform to a well defined notion of normal behavior. Anomaly detection in business process logs has several applications such as fraud detection, identification of

malicious activity and breakdown of the system and understanding the causes of process errors. Due to complex and numerous business processes in a large organizations, it is difficult for any employee to monitor the whole system. As a consequence of this anomalies occurring in a system remains undetected until serious losses are caused by it. Therefore, Root Cause Analysis (RCA) is done to identify root causes and sources of problems and improve or correct the given process so that major problems can be avoided in future.

The focus of the study presented in this thesis is on anomaly detection in business process logs and identification of their root causes. We present a different and fresh perspective to the stated problem and our work is motivated by the need to extend the state-of-the-art in the field of techniques for anomaly detection and RCA in business process event logs. While there has been work done in the area of anomaly detection and RCA in business process logs, to the best of our knowledge, the work presented in this thesis is the first focused study on such a dataset for the application of anomaly detection and RCA. The research aim of the work presented in this thesis is the following:

1. To investigate Window based and Markovian based techniques for detecting anomalies in business process event logs.
2. To apply machine learning techniques such as decision tree classifier to extract rules describing cause of anomalous behavior.
3. To interactively explore different patterns of data using advanced visualization techniques such as parallel plot.
4. To investigate solutions assisting a process analyst to analyze decision tree and parallel plot results, thus identifying root cause of anomalous incidents.
5. To demonstrate the effectiveness of our proposed approach using triangulation study¹. We conduct experiments on a recent, large and real-world incident management data of an enterprise.

¹[http://en.wikipedia.org/wiki/Triangulation_\(social_science\)](http://en.wikipedia.org/wiki/Triangulation_(social_science))

2

Related Work and Novel Research Contributions

In this Chapter, we review work that is closely related to the work presented in this thesis, and list the novel contributions of our work in context to existing work.

2.1 Related Work

Calderón-Ruiz et al. propose a novel technique to identify potential causes of failures in business process by extending available Process Mining techniques [6]. Initially, they filtered the original event log in two logs, the former with successful cases and the latter with failed cases. Then, they extracted behavioral patterns from both event logs using the Performance Sequence Diagram Analysis algorithm [7]. Finally, both sets of patterns are compared considering control flow and time perspectives. The differences found represent potential causes of failures in business processes. They test their technique using several synthetic event logs and are able to successfully find missing or unnecessary activities, and failed behavioural patterns that differ from successful patterns either in the control flow or in the time perspective [6]. Heravizadeh et al. propose a conceptual methodology of root-cause analysis in business processes, based on the definition of softgoals (nonfunctional requirements) for all process activities, as well as correlations between these softgoals and related quality metrics [8]. This methodology requires much effort from participants to document requirements, relationships and respective metrics. Suriadi et al. propose an approach to enrich and transform

process-based logs for Root Cause Analysis based on classification algorithms [9]. They start with determining relevant information that is needed to explain the root cause of a risk incident, followed by the enrichment of the related event log with the necessary information to ensure that sufficient information for RCA is captured [9]. Through the application of aggregation functions and other refinement procedures, they transform enriched event log into a form that is suitable to be analysed by classification techniques. They use decision trees to identify the causes of overtime faults [9]. They validated the applicability of their approach using both self-generated synthetic log and publicly available log [9]. Rogge-Solti et al. propose a Bayesian model that can be automatically inferred from the PetriNet representation of a business process and is then used to detect non-obvious and temporal anomalies [10]. Vasilyev et al. develop an approach to find the cause of delays based on the information recorded in an event log [11]. The approach is based on a logic representation of the event log and on the application of decision tree induction to separate process instances according to their duration [11]. They use inductive logic programming [12], specially a decision tree learner known as TILDE [13] to classify process instances into different groups according to their duration. Each path in decision tree from the root to leaf provides a rule that characterizes a certain group of instances and therefore provides a possible explanation for the delay [11].

Bezerra et al. present some approaches based on incremental mining [14] for anomaly detection, but these algorithms cannot deal with longer traces and/or logs with various classes of traces [4]. Then, in order to deal with such constraints, they begin to develop other solutions based on process mining algorithms available in ProM¹ framework [15] [16]. Bezerra et al. propose an anomaly detection model based on the discovery of an “appropriate process model” [16]. Bezerra et al. apply the process discovery and conformance algorithms from ProM framework for implementing the anomaly detection algorithms [15]. Bezerra et al. present three new algorithms (threshold, iterative, and sampling) to detect “hard to find” anomalies in a process log based only on the control-flow perspective of the traces [17]. This work does not deal with anomalous executions of processes that follow a correct execution path but deal with unusual data, or are executed by unusual roles or users, or have unusual timings [17]. Bezerra et al. develop an algorithm more efficient than the Sampling Algorithm

¹ProM is a pluggable and open-source framework for Process Mining

[18]. They propose an approach for anomaly detection which is an extension of the Threshold Algorithm also reported in [15] [17], which uses process mining tools for process discovery and process analysis for supporting the detection [18].

2.2 Novel Research Contributions

In context to existing work, the study presented in this thesis makes the following novel contributions:

1. Detection of anomalous traces in business process event-logs using Window-based and Markovian-based techniques (after transforming the event-log into a sequential) dataset.
2. Root Cause Analysis (RCA) of anomalous traces using parallel coordinate plots. Application of parallel coordinate plots for visualizing the characteristics of anomalous and normal traces (representing the traces and their attribute values as a polyline with vertices on the parallel axes).
3. Application of tree diagrams as a visual and analytical decision support tool for identifying the features of anomalous traces, thereby assisting a process analyst in problem solving and Root Cause Analysis (RCA).
4. An in-depth and focused empirical analysis on a real-world dataset (Rabobank Group¹: Activity log for incidents) demonstrating the effectiveness of the proposed approach. Application of triangulation technique to validate the outcome of RCA through cross-verification.

¹<http://data.3tu.nl/repository/uuid:c3e5d162-0cfd-4bb0-bd82-af5268819c35>

3

Research Framework and Solution Approach

Figure 3.1 shows the high-level architecture diagram of the proposed solution approach (called as *Pariket*). The proposed approach is a multi-step process primarily consists of 6 phases: experimental dataset collection, sequential dataset conversion, anomaly detection, data pre-processing, classification and visualization. The six phases are labeled in the architecture diagram in Figure 3.1. In phase 1, we download large real world data from Rabobank Group Information and Communication Technology (ICT) (refer to Section 4 on experimental dataset). The dataset is provided in the CSV format. The dataset consists of event logs from interactions records, incidents records, incident activities and change records. The attributes of original CSV files are converted to appropriate data types, such as standardized timestamp formats, for analysis. After loading the data on to MySQL database, we build four tables: Interaction_detail, Incident_detail, Incident_activity_detail and Change_detail. We choose incidents from incident activities to find out anomalous incident patterns. In phase 2, for a particular incident we order the type of activities according to increasing order of DateTime Stamp. Each incident consisting of several activities is represented as a sequence of symbols (refer to Section 5.1 on experimental results). Each unique activity is mapped to a integer symbol. There are 39 different kinds of activities in the dataset and hence there are 39 different symbols. Some of the example of activities are: Referred (28), Problem Closure (22), OOResponse (18), Dial-In (10) and Contact Change (8). The sequences are of different length. The incidents with their corresponding sequence of

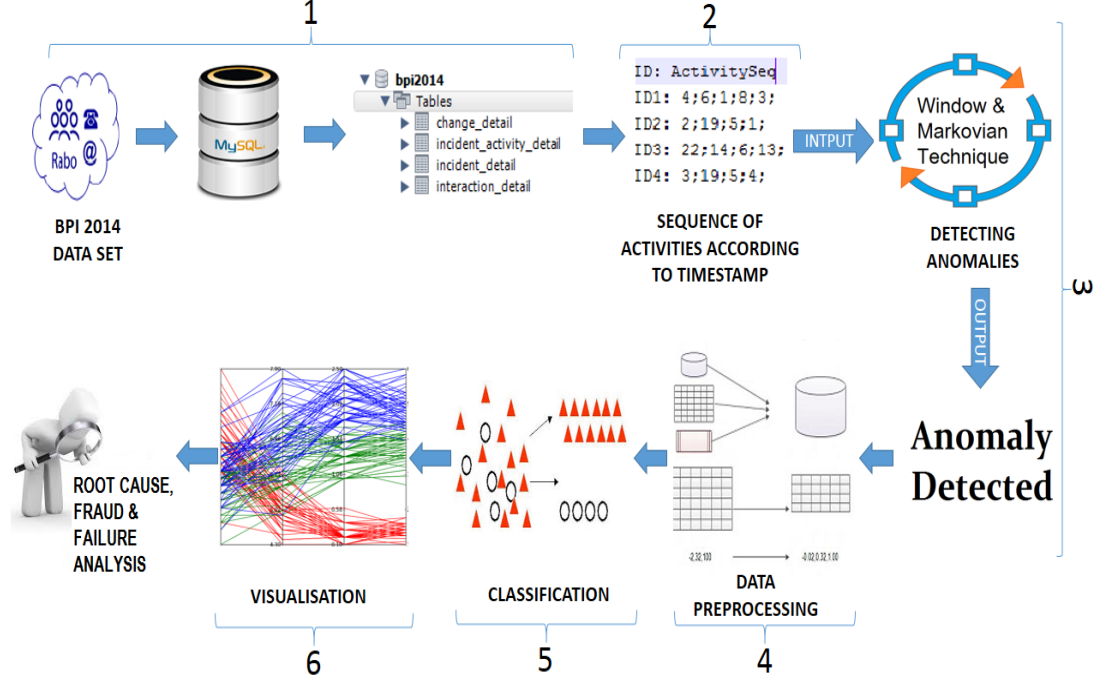


Figure 3.1: Architecture diagram and data processing pipeline for Pariket (Mining Business Process Logs for Root Cause Analysis of Anomalous Incidents)

activities serve as input to anomaly detection algorithms described in Section 5.2.

In phase 3, we implement Window based and Markovian based technique based on the obtained discrete sequences [5] to detect anomalous incidents (refer to Section 5.2 on experimental results). We receive top N anomalous incidents as output from anomaly detection algorithms. We apply decision tree classifier and visualization techniques to identify root causes of anomalous incidents. Input to these techniques requires data to be in a particular format and of high quality. Hence, in phase 4 we perform data pre-processing to bring the data in the required format and of high quality. The data pre-processing helps in improving the accuracy and efficiency of the subsequent mining processes. Data goes through series of steps during pre-processing phase: integration, cleaning and transformation etc (refer to Section 5.3 of experimental results). In phase 5, we create decision tree using $J48$ algorithm in Waikato Environment for Knowledge Analysis (Weka)¹ to identify the features of anomalous traces, thereby assisting a pro-

¹<http://www.cs.waikato.ac.nz/ml/weka/>

cess analyst in problem solving and Root Cause Analysis (RCA) (refer to Section 5.4). We choose Weka as it provides a flexible interface which is easy to use. The *J48* tree classifier is the *C4.5* implementation available in Weka. *J48* handles both numeric and nominal attribute values. In phase 6, we apply advanced visualization technique such as parallel plot in Tibco Spotfire¹ to interactively explore characteristics of anomalous and normal traces (refer to Section 5.5). Tibco Spotfire is an analytics software that helps quickly uncover sights for better decision making. Business process analyst then analyze decision tree and parallel plot results to identify the root cause of anomalous incidents.

¹<http://spotfire.tibco.com/>

4

Experimental Dataset

We conduct our study on a large real world data from Rabobank Group Information and Communication Technology (ICT). The data is related to the Information Technology Infrastructure Library (ITIL) process implemented in the Bank. The ITIL process depicted in Figure 4.1 starts when an internal client reports an issue regarding disruption of ICT service to Service Desk Agent (SDA). SDA records the complete information about the problem in Interaction Record. If the issue does not get resolved on first contact then a Incident record is created for the corresponding Interaction else the issue is closed. There can be many to one mapping between Interaction Record and Incident Record. If a issue appears frequently then a request for change is initiated.

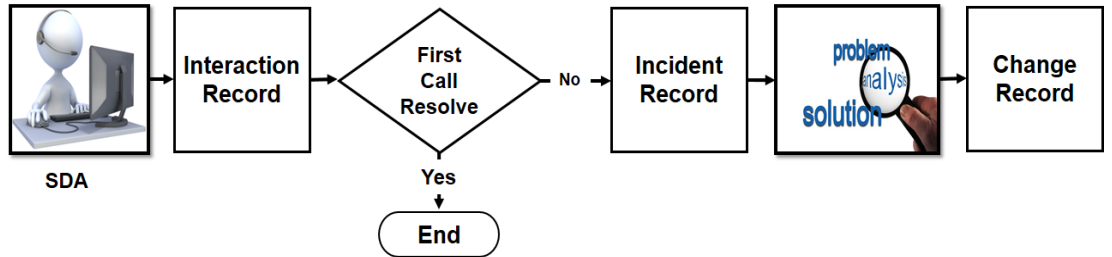


Figure 4.1: Information Technology Infrastructure Library (ITIL) process implemented in Rabobank Group

The dataset is provided in the CSV format. It contains the event logs from interactions records, incidents records, incident activities and change records. The provided dataset is of six month duration from October 2013 - March 2014. Interactions that were not resolved before 31 March, were removed from the dataset. The attributes of

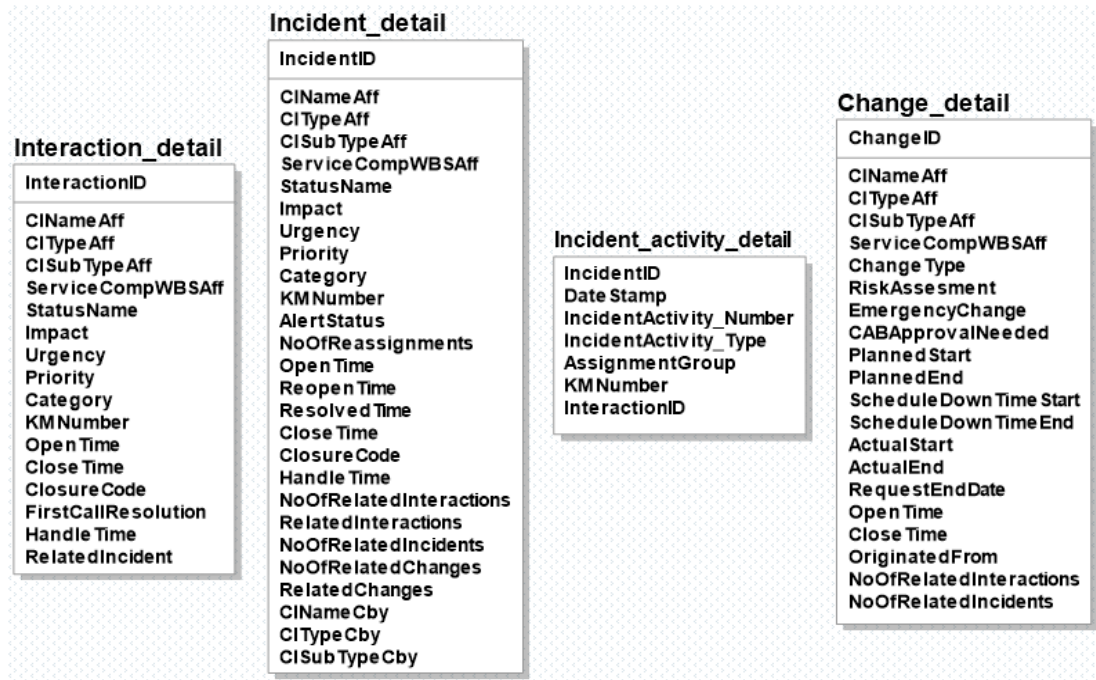


Figure 4.2: Output of CA ERwin Data Modeler

original .CSV files are converted to appropriate data types, such as standardized timestamp formats, for analysis. After loading the data on to MySQL database, we build four tables: Interaction_detail, Incident_detail, Incident_activity_detail and Change_detail.

Figure 4.2 is the representation of tables in our database which is created using CA ERwin Data Modeler tool¹.

1. Interaction_detail - It has 147,004 records, each one corresponding to an interaction. Every record contains information like InteractionID, Priority, Category, Open Time, Close Time, Handle Time, and First Call Resolution (whether SDA was able to resolve the issue on first contact or not).
2. Incident_detail - It has 46,606 records, each one corresponding to an incident case. Every record has attributes like IncidentID, Related Interaction, Priority, Open Time, Handle Time, Configuration Item Affected etc.
3. Incident_activity_detail - It has 466,737 records. Each record contains an IncidentID with the activities performed on it. It also contains information about

¹ <http://erwin.com/products/datamodeler>

the Assignment Group that is responsible for a particular activity.

4. Change_detail - It contains records of the activities performed on each change case. It has information about Configuration Item Affected, Service Component Affected, Change Type and Risk Assessment etc.

As an academic, we believe and encourage academic code or software sharing in the interest of improving *openness and research reproducibility*. We release our code and dataset in public domain so that other researchers can validate our scientific claims and use our tool for comparison or benchmarking purposes (and also reusability and extension). Our code and dataset is hosted on GitHub¹ which is a popular web-based hosting service for software development projects. We select GPL license (restrictive license) so that our code can never be closed-sourced.

¹<https://github.com/ashishsureka/pariket>

5

Experimental Results

We perform a series of steps to identify the root cause of anomalous incidents. Each of the following 6 sub-sections describes the steps consisting of procedure or approach and findings.

5.1 Sequential Dataset Conversion

We analyze all the tables and amongst them we choose Incident_activity_detail table (refer Figure 4.2 in Section 4) to find out the anomalous incident patterns. This table contains a log of activities performed by the service team(s) to resolve incidents which are not resolved by first contact. The main reason for choosing this table is because it has information regarding the type of activities performed on a particular incident id and also the timestamp when this incident activity type started to resolve the issue.

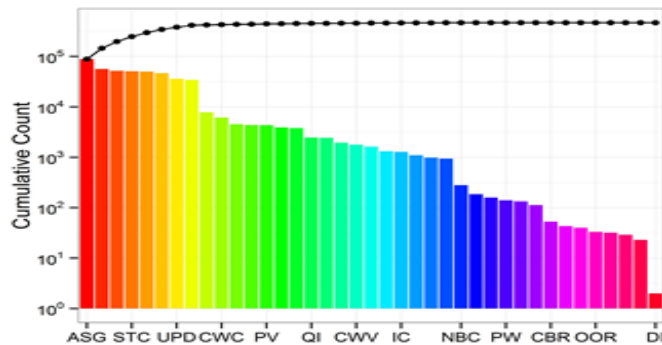


Figure 5.1: Pareto chart showing the distribution of activities and their cumulative count. Y-axis is in logarithmic scale

5.1 Sequential Dataset Conversion

IncidentID	DateStamp	IncidentActivity_Number	IncidentActivity_Type	AssignmentGroup	KMnumber	InteractionID	IncidentActivity_Type_Number
IM0000004	2013-01-07 08:17:17	001A3689763	Reassignment	TEAM0001	KM0000553	SD0000007	27
IM0000004	2013-11-04 13:41:30	001A5852941	Reassignment	TEAM0002	KM0000553	SD0000007	27
IM0000004	2013-11-04 13:41:30	001A5852943	Update from customer	TEAM0002	KM0000553	SD0000007	34
IM0000004	2013-11-04 12:09:37	001A5849980	Operator Update	TEAM0003	KM0000553	SD0000007	20
IM0000004	2013-11-04 12:09:37	001A5849979	Assignment	TEAM0003	KM0000553	SD0000007	2
IM0000004	2013-11-04 13:41:30	001A5852942	Assignment	TEAM0002	KM0000553	SD0000007	2
IM0000004	2013-11-04 13:51:18	001A5852172	Closed	TEAM0003	KM0000553	SD0000007	5
IM0000004	2013-11-04 13:51:18	001A5852173	Caused By CI	TEAM0003	KM0000553	SD0000007	4
IM0000004	2013-11-04 12:09:37	001A5849978	Reassignment	TEAM0003	KM0000553	SD0000007	27
IM0000004	2013-09-25 08:27:40	001A5544096	Operator Update	TEAM0003	KM0000553	SD0000007	20

↓
Order activities according to date timestamp
↓

IncidentID	DateStamp	IncidentActivity_Number	IncidentActivity_Type	AssignmentGroup	KMnumber	InteractionID	IncidentActivity_Type_Number
IM0000004	2013-01-07 08:17:17	001A3689763	Reassignment	TEAM0001	KM0000553	SD0000007	27
IM0000004	2013-09-25 08:27:40	001A5544096	Operator Update	TEAM0003	KM0000553	SD0000007	20
IM0000004	2013-11-04 12:09:37	001A5849980	Operator Update	TEAM0003	KM0000553	SD0000007	20
IM0000004	2013-11-04 12:09:37	001A5849979	Assignment	TEAM0003	KM0000553	SD0000007	2
IM0000004	2013-11-04 12:09:37	001A5849978	Reassignment	TEAM0003	KM0000553	SD0000007	27
IM0000004	2013-11-04 13:41:30	001A5852941	Reassignment	TEAM0002	KM0000553	SD0000007	27
IM0000004	2013-11-04 13:41:30	001A5852943	Update from customer	TEAM0002	KM0000553	SD0000007	34
IM0000004	2013-11-04 13:41:30	001A5852942	Assignment	TEAM0002	KM0000553	SD0000007	2
IM0000004	2013-11-04 13:51:18	001A5852172	Closed	TEAM0003	KM0000553	SD0000007	5
IM0000004	2013-11-04 13:51:18	001A5852173	Caused By CI	TEAM0003	KM0000553	SD0000007	4

Figure 5.2: Activities ordered according to increasing order of datetime stamp for incident id IM0000004

The attribute IncidentActivity_Type represents the type of activity performed on the incident. There are 39 unique activities. Some of the examples are: Assignment (ASG), Status Change (STC), Update (UPD), Referred (REF), Problem Closure (PC), OOResponse (OOR), Dial-In (DI) and Contact Change (CC). Figure 5.1 represents the pareto chart showing the distribution of activities and their cumulative count. The Y-axis is in logarithmic scale. We assign integer number starting with 0 to 38 to these activities, and then we add an extra column IncidentActivity_Type_Number into the table Incident_activity_detail denoting this activity number. For a particular incident we order the activities according to increasing order of DateTime Stamp. This is done for all the unique incidents in the Incident_activity_detail table. Figure 5.2 shows the screenshot from MySql of activities performed during one of the incidents ‘IM0000004’.

The even-log data in Figure 5.2 shows that several activities are performed by vari-

IncidentID	IncidentActivity_Type_List
IM0000004	27;20;20;2;27;27;34;2;5;4;
IM0000005	27;20;9;27;20;2;20;32;32;2;32;27;1;2;2;33;2;20;27;2;20;34;27;2;34;32;4;5;
IM0000011	27;2;27;20;34;2;34;27;2;34;2;20;27;34;34;2;27;2;5;4;
IM0000012	34;27;2;34;4;5;

Figure 5.3: Incident id's with the sequence of activities separated by semicolon

ous actors during the work-flow and process enactment. Figure 5.2 shows that the data has a sequential aspect (is a nature and characteristics of the business process log) and hence we believe techniques for anomaly detection for sequences can be applied to the event log data. While the sequence in the given example is multivariate, in this work, we consider only the activity attribute and model the sequence as univariate. Figure 5.3 shows screenshot from MySQL, each incident consisting of several activities is represented as a sequence of integer numbers. We apply Window Based and Markovian Based Techniques for detecting anomalous incidents. The input dataset to these algorithms has to be in sequential format. Therefore, to accomplish this we create a new table 'IncidentActivitySequence' containing two attributes 'IncidentID' and 'IncidentActivity_Type_List'. Each record in this table contains all unique IncidentID's and sequences of IncidentActivity_Type_Number separated by semicolon according to timestamp from Incident_activity_detail. For example, corresponding to IncidentID 'IM0000012', the sequence of activities are '34;27;2;34;4;5;'.

5.2 Anomaly Detection

The outcome of the Section 5.1 is a list of all incident id's with their corresponding sequence of activities ordered according to timestamp. The aim of algorithms described in this Subsection is to identify anomalous incidents based on the obtained discrete sequences. There is no reference or training database available containing only normal sequences. Hence, our task is to detect anomalous sequences from an unlabeled database of sequences. The problem is of unsupervised anomaly detection. A formal representation of the problem is [5]: Given a set of n sequences, $S = \{S_1, S_2, \dots, S_n\}$, find all sequences in S that are anomalous with respect to rest of S . This unsupervised

problem can be solved by using a semi-supervised approach where we treat the entire dataset as training set and then score each sequence with respect to this training set. We assume that majority of sequences in the unlabeled database are normal as anomalies are generally infrequent in nature [5]. We use two algorithms, Window Based and Markovian Based described in the following Subsections for anomaly detection.

5.2.1 Window Based Technique

The motivation behind using window based technique is to determine anomalous sequences even if the cause of anomaly is localized to one or more shorter subsequences within the actual sequence [19]. Window based technique in general operates as, first we extract overlapping windows of fixed length (k) from a given test sequence. Then, we assign some anomaly score to each extracted window based on a threshold value (λ). Finally, the anomaly score of all the windows are combined to obtain an anomaly score for the test sequence [5].

The pseudocode for Window Based anomaly detection algorithm is shown in Algorithm 1. The input to the algorithm is data comprising of IncidentID, sequence of activities from table IncidentActivitySequence, window size (k), threshold (λ) and number of anomalous incidents (N). The algorithm returns top N anomalous IncidentID as output. The main challenge was to find out the size of window (k) and the value of threshold (λ). We analyze all the subsequences of window length less than 3. Our analysis reveals that they occur very frequently. Therefore, we cannot take them as anomalous subsequences because according to our previous assumption in Section 5.2 anomalies in our dataset are in minority. Therefore, k has to be equal to or greater than 3.

Algorithm 1 consists of two phases: training and testing. We choose 3 experimental parameters: $k = 3$, $\lambda = 4$ and $N = 1000$. The training phase is represented by Steps 4-12. During this phase, we obtain the sequence of activities for each IncidentID. From the sequence we extract k length overlapping (sliding) windows. We maintain each unique window with its frequency in normal dictionary D . The testing phase is represented by Steps 13-25. Every sequence of the training dataset is considered as the test sequence. During this phase, we extract sliding windows of length k from the test sequence S_i . A window W_j is assigned an anomalyScore of 1 if the frequency associated with the window W_j in dictionary D is less than the threshold value (λ) else anomalyScore is 0.

Algorithm 1: Window Based Algorithm (ID, S, k , λ , N)

Data: IncidentID ($ID = ID_1 \dots ID_n$) and Sequence of activities ($S = S_1 \dots S_n$) from table.

Result: Top N Anomalous IncidentID.

```
1 set windowSize =  $k$ , threshold =  $\lambda$ ;  
2 create a empty dictionary  $D$ ,  $D'$   
3 create an arrayList anomalousIncidents;  
4 foreach IncidentID  $ID_i$  in ID do  
5    $S_i$  = get the sequence corresponding to  $ID_i$   
6   set windowCount =  $S_i$ .length - windowSize + 1  
7   foreach  $j = 1$  to windowCount do  
8     read the subsequence ( $W_j$ ) of length = windowSize starting from  $j^{th}$   
       position in  $S_i$   
9     if  $W_j$  is not present in  $D$  then  
10       $\lfloor$  add ( $W_j$ , 1) as (key, value) pair in  $D$   
11    else  
12       $\lfloor$  add ( $W_j$ , value + 1) in  $D$   
13 foreach IncidentID  $ID_i$  in ID do  
14    $S_i$  = get the sequence corresponding to  $ID_i$   
15   set windowCount =  $S_i$ .length - windowSize + 1  
16   set anomalyScore = 0.0  
17   foreach  $j = 1$  to windowCount do  
18     read the subsequence ( $W_j$ ) of length = windowSize starting from  $j^{th}$   
       position in  $S_i$   
19     get the (key, value) pair from  $D$  corresponding to key =  $W_j$   
20     if value is less than threshold then  
21        $\lfloor$  anomalyScore = anomalyScore + 1  
22   anomalyScore = anomalyScore / windowCount  
23   add  $ID_i$ , anomalyScore into  $D'$   
24 sort  $D'$  according to decreasing anomalyScore  
25 add top  $N$  IncidentId into anomalousIncidents  
26 return anomalousIncidents
```

To calculate the anomalyScore of a complete test sequence S_i , we take summation of anomalyScore of all the subsequence windows contained in it. The anomaly score of the test sequence is proportional to the number of anomalous windows in the test sequence [20]. The result obtained after executing Steps 14-21 is then divided by the number of windows contained in the test sequence. This normalization is done to take into account the varying lengths of sequences. The anomalyScore of 1 for a test sequence denotes most anomalous and 0 as least anomalous. We store the IncidentID and its corresponding anomalyScore in the dictionary D' . Then, we sort the IncidentID's in decreasing order of anomalyScore and return top N anomalous IncidentID's.

5.2.2 Markovian Based Technique

We apply fixed Markovian technique [5] which is based on the property of short memory of sequences. This property states that the conditional probability of occurrence of a symbol s_i is dependent on the occurrence of previous k symbols with in a sequence S_i [21]. The conditional probability of occurrence of a symbol s_i in a sequence S_i is given by Equation 5.1:

$$P(s_i | s_{(i-k)} \dots s_{(i-1)}) = \frac{\text{freq}(s_{(i-k)} \dots s_i)}{\text{freq}(s_{(i-k)} \dots s_{(i-1)})} \quad (5.1)$$

where $\text{freq}(s_{(i-k)} \dots s_i)$ is the frequency of occurrence of the subsequence $s_{(i-k)} \dots s_i$ in the sequences in S and $\text{freq}(s_{(i-k)} \dots s_{(i-1)})$ is the frequency of occurrence of the subsequence $s_{(i-k)} \dots s_{(i-1)}$ in the sequences in S.

The pseudocode for Markovian based anomaly detection algorithm is shown in Algorithm 2. The input to the algorithm is data comprising of IncidentID, sequence of activities from table IncidentActivitySequence, window size (k) and number of anomalous incidents (N). The algorithm returns top N anomalous IncidentID as output. Algorithm 2 consists of two phases: training and testing. Steps 4-17 represents the training phase. During this phase, we create two dictionaries D_k and D_{k+1} of length k and $k+1$ respectively. The process for creation of dictionary is similar to that described in Section 5.2.1. We choose $k = 3$ for our experiment. It takes into account the subsequences of length 4 which are dependent on previous 3 symbols. Steps 18-32 represents the testing phase. Steps 18-32 are repeated for each IncidentID in table IncidentActivitySequence. We extract the test sequence S_i corresponding to a IncidentID ID_i in

Algorithm 2: Markovian Based Algorithm (ID, S, k , N)

Data: IncidentID (ID= $ID_1 \dots ID_n$) and Sequence of activities (S= $S_1 \dots S_n$)
from table.

Result: Top N Anomalous IncidentID.

```

1 create a empty dictionary  $D_k, D_{k+1}, D'$ 
2 create an arrayList anomalousIncidents
3 foreach IncidentID  $ID_i$  in ID do
4    $S_i$  = get the sequence corresponding to  $ID_i$ 
5   set noOfSubsequences =  $S_i.length - k + 1$ 
6   foreach  $j = 1$  to noOfSubsequences do
7     read the subsequence ( $W_j$ ) of length =  $k$  starting from  $j^{th}$  position in  $S_i$ 
8     read the subsequence ( $W_{j+1}$ ) of length =  $k+1$  starting from  $j^{th}$  position
      in  $S_i$ 
9     if  $W_j$  is not present in  $D_k$  then
10      | add ( $W_j, 1$ ) as (key, value) pair in  $D_k$ 
11    else
12      | add ( $W_j, value + 1$ ) in  $D_k$ 
13    if  $W_{j+1}$  is not present in  $D_{k+1}$  then
14      | add ( $W_{j+1}, 1$ ) as (key, value) pair in  $D_{k+1}$ 
15    else
16      | add ( $W_{j+1}, value + 1$ ) in  $D_{k+1}$ 
17 foreach IncidentID  $ID_i$  in ID do
18    $S_i$  = get the sequence corresponding to  $ID_i$ 
19   set noOfSubsequences =  $S_i.length - k + 1$ 
20   set anomalyScore = 0.0, prob = 0
21   foreach  $j = 1$  to noOfSubsequences - 1 do
22     read the subsequence ( $W_j$ ) of length =  $k$  starting from  $j^{th}$  position in  $S_i$ 
23     read the subsequence ( $W_{j+1}$ ) of length =  $k+1$  starting from  $j^{th}$  position
      in  $S_i$ 
24     get the ( $key_j, value_j$ ) pair from  $D_k$  corresponding to key =  $W_j$ 
25     get the ( $key_{j+1}, value_{j+1}$ ) pair from  $D_{k+1}$  corresponding to key =  $W_{j+1}$ 
26      $r = (value_j) / (value_{j+1})$ ;
27     prob = prob + log( $r$ );
28   prob = prob / noOfSubsequences
29   TestSequenceProbability =  $e^{prob}$ 
30   anomalyScore = 1 / TestSequenceProbability;
31   add  $ID_i$ , anomalyScore into  $D'$ 
32 sort  $D'$  according to decreasing anomalyScore
33 add top  $N$  IncidentID into anomalousIncidents
34 return anomalousIncidents

```

Step 19 and calculate the number of subsequences of length k in Step 20. Steps 23-28 are repeated for each each subsequence within the test sequence S_i . Step 23 and 24 reads the subsequence W_j and W_{j+1} of length k and $K+1$ respectively starting from position j . Step 25 and 26 calculates the frequency $value_j$ and $value_{j+1}$ of W_j and W_{j+1} from the dictionaries D_k and D_{k+1} . We calculate the conditional probabilities of symbols in Step 27 by using Equation 5.1. We calculate the overall probability of S_i using the Equation 5.2:

$$P(S_i) = \prod_{i=1}^l P(s_i | s_1 s_2 \dots s_{i-1}) \quad (5.2)$$

where l is the length of the sequence S_i and s_i is the symbol occurring at position i in S_i [5]. For simplification, we take \log on both the sides in Equation 5.2, the modified equation is used in the Step 28. We normalize the probability in Step 29 to take into account the varying length of sequences. We calculate anomaly score for test sequence S_i as the inverse of the probability of S_i in Step 31. Less probability of the test sequence means more anomaly score. We store the IncidentID and its corresponding anomalyScore in the dictionary D' . Then, we sort the IncidentID's in decreasing order of anomalyScore and return top N IncidentID's.

5.3 Data Pre-processing

We receive top N anomalous IncidentID's as output from algorithms described in Section 5.2. Our aim is to identify root causes of anomalous incidents. We apply data mining techniques (machine learning) such as decision tree classifier to extract rules describing anomalous behaviour. Input to these techniques requires data to be in a particular format and of high quality. Hence, we apply data pre-processing techniques to bring the data in the required format and also to improve the quality. The data pre-processing helps in improving the accuracy and efficiency of the subsequent mining processes. Data goes through series of steps during pre-processing phase: integration, cleaning and transformation etc.

First, we join two tables 'Interaction_detail' and 'Incident_detail' (refer to Phase 4 of architecture diagram in Figure 3.1). We use attribute 'RelatedIncident' from Interaction_detail table as foreign key and 'IncidentID' from Incident_detail as primary key to perform the join. We give new name 'Interaction_Incident' to the merged

Table 5.1: Name, Type and Description of Some of Attributes in Merged table of Interaction_detail and Incident_detail

Attribute Name	Attribute Type	Description
Incident_CIType(Aff)	Nominal	There are 13 distinct types of CIs. Example: software, storage, database, hardware, application.
Incident_CISubType(Aff)	Nominal	There are 64 CI Sub-types. Example : web based, client based, server based, SAP.
Incident_Priority	Nominal	There are 5 categories of priority i.e {1, 2, 3, 4, 5}.
Incident_Category	Nominal	There are 4 Incident Category i.e {Incident, Request For Information, Complaint, Request For Change}
Incident_OpenTime	Date Format	The Open Time of Incident is in 'yyyy-MM-dd HH:mm:ss' format. We convert into timestamp in 'hours'.
Incident_HandleTime	Date Format	The Handle Time of Incident is in 'yyyy-MM-dd HH:mm:ss' format. We convert into timestamp in 'hours'.
Interaction_CIType(Cby)	Nominal	There are 13 distinct types of caused by CIs. Example: software, storage, database, hardware, application.
Interaction_CISubType(Cby)	Nominal	There are 64 CI Sub-types. Example: web based, client based, server based, SAP.
Closure Code	Nominal	There 15 distinct types of Closure Code. Example: Unknown, Operator Error, Enquiry, Hardware, Software.
Anomalous	Nominal	{Yes, No}.

table. The merged table contains all the information for the issues that could not be resolved on first call. We create two copies of table `Interaction_Incident`: `Interaction_Incident_Markovian` and `Interaction_Incident_Windows`. We add new attribute ‘Anomalous’ to the newly created tables. We make the value of attribute ‘Anomalous’ as ‘Yes’ for all the top N anomalous IncidentID’s and ‘No’ for rest of the records. Table 5.1 represents name, type and description of some of attributes obtained after merging `Interaction_detail` and `Incident_detail`. The anomalous IncidentID’s for table `Interaction_Incident_Windows` are obtained from the outcome of Algorithm 1. The anomalous IncidentID’s for table `Interaction_Incident_Markovian` are obtained from the outcome of Algorithm 2. We use J48 algorithm for classification using decision tree in Weka. The J48 algorithm handles missing values itself by replacing them with the most frequent observed non-missing values.

Next we transform open time, close time for `Interaction` and open time, reopen time, resolved time, closed time for `Incident` given in datetime format. We covert the datetime format that is ‘yyyy-MM-dd HH:mm:ss’ into timestamp in ‘hours’ to be useful for classification. For this, we take reference datetime as ‘1970-01-01 17:13:01’. Handle time for both `Interaction` and `Incident` is given in seconds in comma separated format. We transform it by removing comma because the J48 algorithm takes input in CSV or Attribute-Relation File Format (ARFF) format. The pre-processed data serves as input to the classification and visualization techniques.

5.4 Classification

Data mining technique such as decision tree offer a semi automated approach to identify root causes of anomalous incidents. Choosing a data mining analysis tool to execute decision tree algorithm can be a challenge. Popular open source data mining packages include Weka, R, Tanagra, Yet Another Learning Environment (YALE), and Konstanz Information Miner (KNIME) [22]. We choose Weka [23] as it provides a flexible interface which is easy to use. Figure 5.4 shows the graphical user interface of Weka. Weka is open source software issued under the GNU General Public License. Weka is a collection of machine learning algorithms for data mining tasks like data pre-processing, classification, clustering, association rules, visualization, etc.

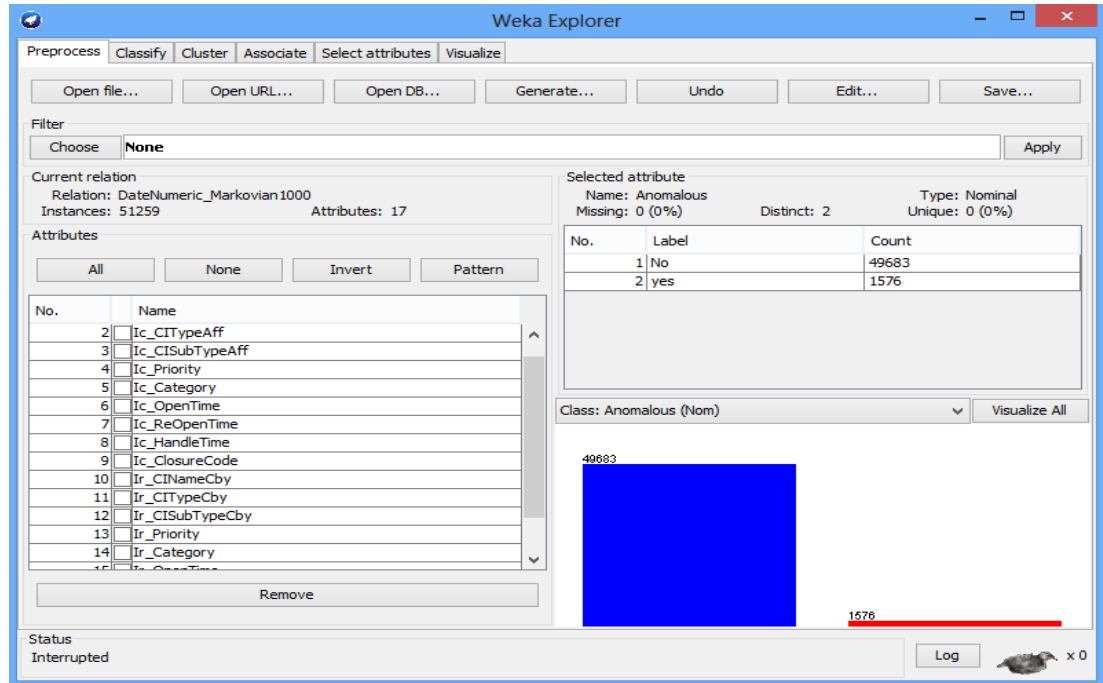


Figure 5.4: Graphical user interface of Weka

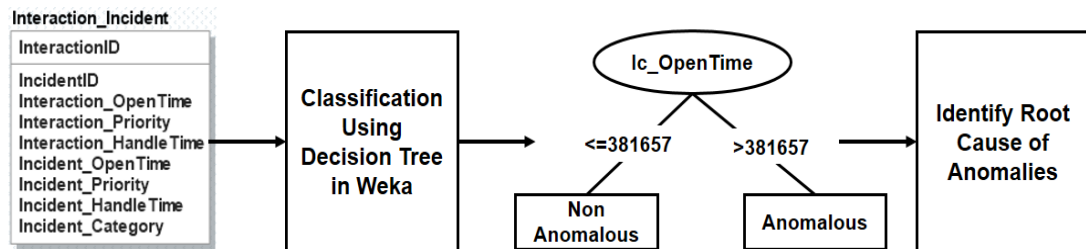


Figure 5.5: Flow of classification using Weka

Figure 5.5 depicts overall flow of classification using decision tree in Weka. The pre-processed data which we obtain after data pre-processing in Section 5.3 serves as input to Weka. The input to Weka is normally in CSV or ARFF format. We use Attribute Selector to select attributes in Weka. Attribute selection involves searching through all possible combination of attributes in the data to find which subset attributes works best for prediction. Attribute selection process is separated into two parts: Attribute Evaluator and Search Method. Attribute Evaluator has methods to assess attribute subsets. Some examples of attribute evaluation methods are CfsSubsetEval, Classifier-

SubsetEval, InfoGainAttributeEval, etc. Search Method is the method by which the space of possible attribute subsets is searched. Some examples of search methods are Random Search, Exhaustive Search, Best First Search, etc. We fed the classification ready pre-processed data into Weka. Weka supports classification algorithms, such as J48 [24], JRip [25], and many others. We perform classification using decision tree algorithm in Weka. Decision tree offers many benefits: easy to understand by user, handles variety of input data such as nominal and numeric and handles missing values in dataset. We perform classification using different algorithms for decision tree like J48, ADTree, REPTree and SimpleCart [26].

Algorithm Used	Markovian Technique Confusion Matrix	Window Based Technique Confusion Matrix
J48	a b < -- classified as 49666 17 a= No 1131 445 b= yes	a b < -- classified as 49943 6 a= No 1146 164 b= yes
ADTree	a b < -- classified as 49673 10 a= No 1266 310 b= yes	a b < -- classified as 49938 11 a= No 1242 68 b= yes
REPTree	a b < -- classified as 49650 33 a= No 1141 435 b= yes	a b < -- classified as 49916 33 a= No 1156 154 b= yes
SimpleCart	a b < -- classified as 49656 27 a= No 1164 412 b= yes	a b < -- classified as 49934 15 a= No 1179 131 b= yes

Figure 5.6: Decision tree classifier algorithm names and confusion matrices on records from Markovian and Window based technique

Figure 5.6 shows the result of different classification algorithms on records obtained from Markovian and Window Based technique. We observe that J48 algorithm has higher true positive and true negative rate in comparison to other algorithms. Therefore, we apply J48 algorithm for decision tree based classification on data. The J48 tree classifier is the C4.5 implementation available in Weka. The J48 builds decision tree from a set of labeled training data using the concept of information entropy. We change the default parameters in J48 algorithm like binarySplits, ConfidenceFactor,

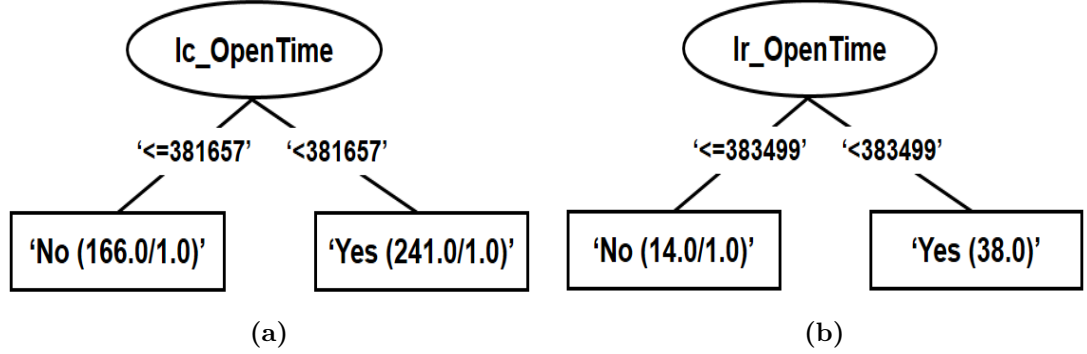


Figure 5.7: Fragments of decision tree using Weka based on anomalous incidents received from Markovian based technique

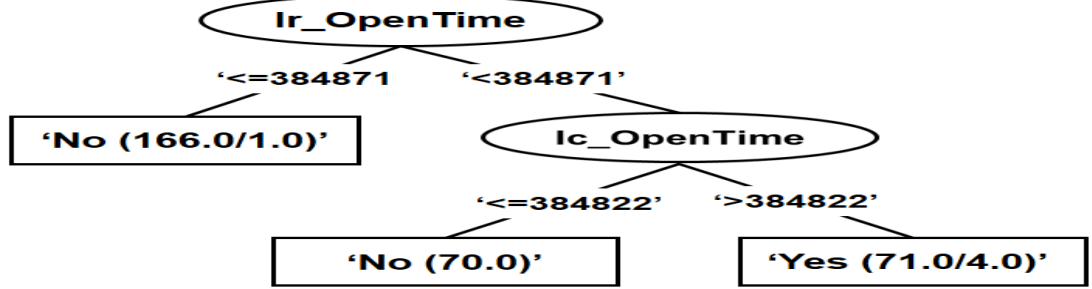


Figure 5.8: Fragment of decision tree using Weka based on anomalous incidents received from Window based technique

minNumObj, etc. But, there is no improvement observed in the results. Result is displayed in classifier output window. To view tree in graphical format click on 'visualize tree' option in pop menu.

We consider combination of attributes or parameters as root causes of anomalous incidents which occur on the path from the root to leaf showing anomalous as *Yes*. Figure 5.9 shows the decision tree using the J48 algorithm under 10-fold cross validation mode on records from Interaction_Incident_Markovian in CSV format. Figure 5.7 and 5.8 shows the fragments of the decision tree extracted from Weka (due to limited space it is not possible to display the entire tree). To represent figures more clearly, Incident is written as Ic and Interaction is written as Ir. We create decision tree in Figure 5.7a and 5.7b by applying J48 algorithm on records from Interaction_Incident_Markovian in CSV format. We observe that there are 240 anomalous incidents whose incident open time is greater than 381657 (hrs). Figure 5.7b shows that there are 38 anomalous

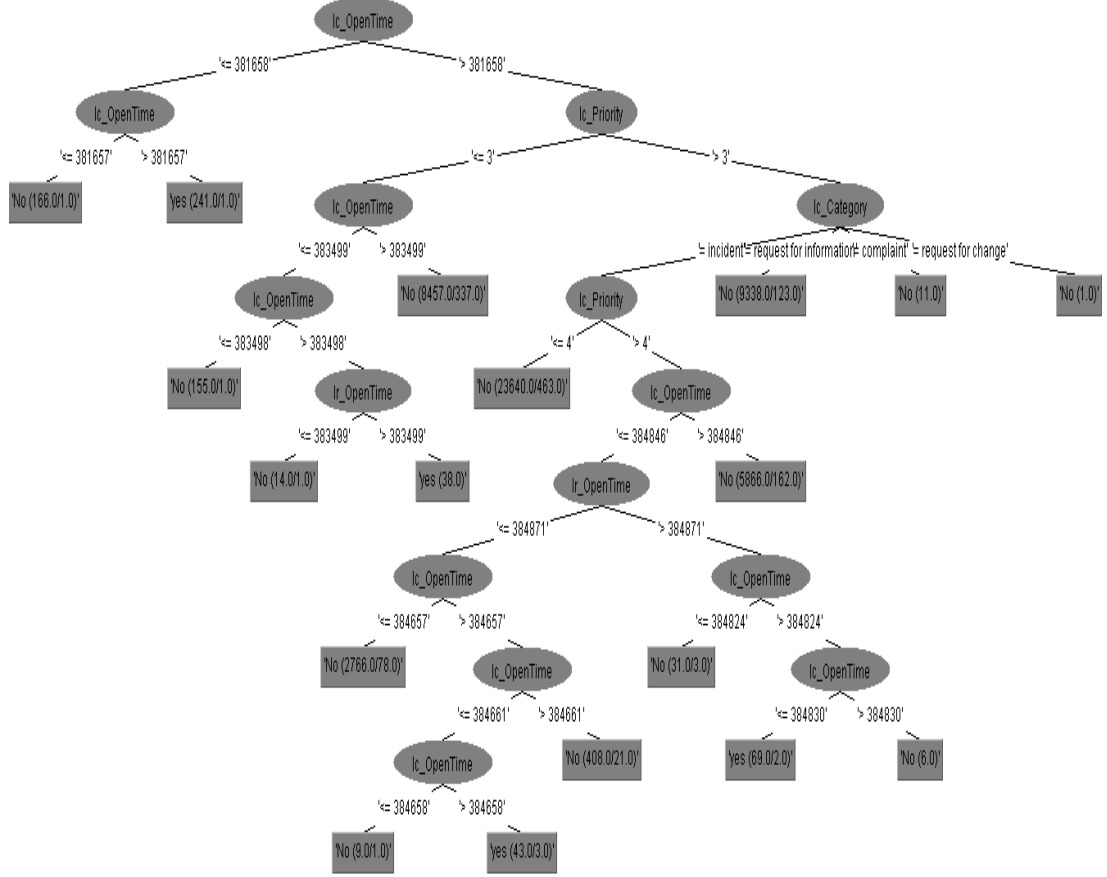


Figure 5.9: Decision tree using Weka based on anomalous incidents received from Markovian based technique

incidents whose interaction open time is greater than 383499 (hrs). Decision tree in Figure 5.8 is for results from Interaction_Incident_Windows in CSV format. Figure 5.8 depicts there are 67 anomalous incidents whose interaction open time is greater than 384871 (hrs) and incident open time is greater than 384822 (hrs).

We obtain attributes Incident_OpenTime, Interaction_OpenTime, Incident_Priority and Incident_Category on the path which leads to anomalous incident leaf nodes. And, there is a path consisting of only Incident_OpenTime which classifies 240 incidents as anomalous. Therefore, open time of incidents alone or combination of open time of interaction, open time of incident, priority of incident and category of incident are causes behind anomalous incidents.

5.5 Visualization

Visualization techniques are used to facilitate user interaction with data. User analyze data by carefully examining it and using different tools on it. Visualization techniques help to identify usual trends and anomalies which are present in data. To achieve this, We use the Tibco's Spotfire platform. Tibco Spotfire is an analytics software that helps quickly uncover sights for better decision making. It is used to detect patterns and correlations present in the data that were hidden in our previous approach using Decision tree. Among many features provided by Spotfire, we use Parallel Coordinate Plot for visualization.

Parallel Coordinate Plot maps each row in the data table as a line. Each attribute of a row is represented by a point on the line. The values in Parallel Plot are normalized. It means lowest value for an attribute in the column is 0% of entire data values in that column while the highest value is 100% unlike the line graphs. Therefore, we cannot compare the values in one column with the values in other column. The data fed into parallel plot is the integrated data which we obtain after the pre-processing phase described in Section 5.3. We create parallel plot by following four steps.

1. Load the data into the Tibco Spotfire. Data can be of type Spotfire Binary Data Format (SBDF), TXT, XLSX, CSV. etc.
2. Choose Parallel Coordinate Plot from Insert tab.
3. Select attributes from column option of the properties section. The selected attributes will be displayed on X-axis of plot.

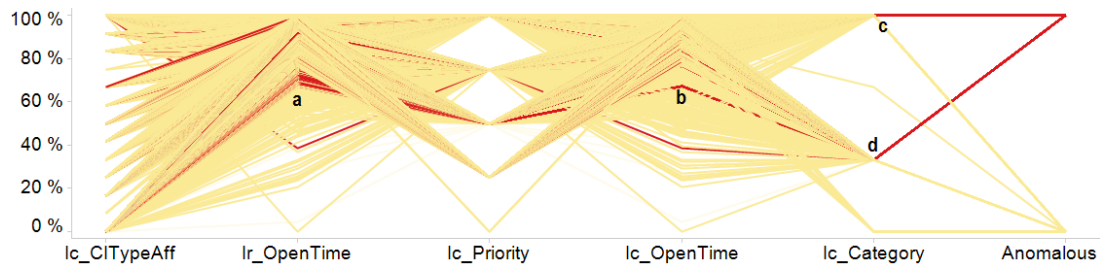


Figure 5.10: Parallel coordinate plot depicting the behaviour of incident CI type affected, interaction open time, incident priority, incident open time and incident category for anomalous and non-anomalous incidents from Markovian technique

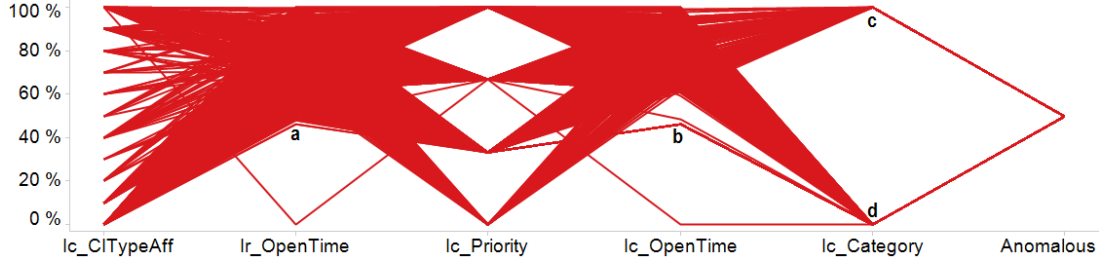


Figure 5.11: Parallel coordinate plot depicting the behaviour of incident CI type affected, interaction open time, incident priority, incident open time and incident category for anomalous incidents from Markovian technique

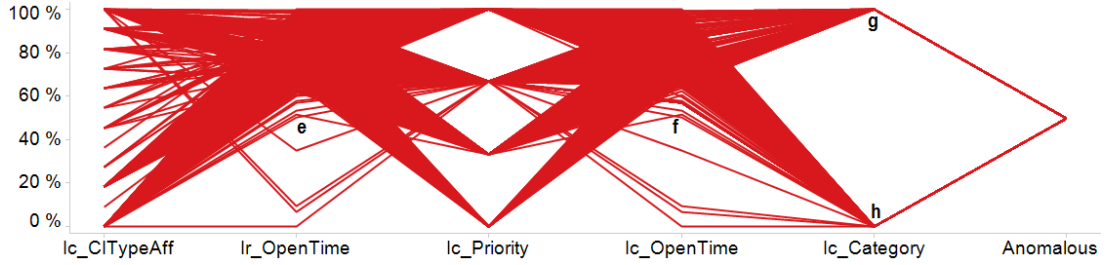


Figure 5.12: Parallel coordinate plot depicting the behaviour of incident CI type affected, interaction open time, incident priority, incident open time and incident category for anomalous incidents from Window based technique

4. Color the lines or profile of each data row depending on attribute value.

We create plots by taking into account different combination of attributes along with the attribute Anomalous (Yes/No). Figure 5.10 represents patterns of the attributes: Ic_CITyPeAff, Ir_OpenTime, Ic_Priority, Ic_OpenTime and Ic_Category for the complete dataset on parallel plot. The dataset consists of records from table Interaction_Incident_Markovian created in Pre-processing phase in CSV format. Due to scarcity of space, Incident is written as Ic and Interaction is written as Ir. We show anomalous incidents with red color and non-anomalous with yellow color. We consider only anomalous incidents in Figure 5.11 for better clarity. The attributes Ic_CITyPeAff and Ic_Priority individually do not show any useful information regarding anomalous behavior of incidents. Anomalous Incidents are falling in all the Ic_CITyPeAffs and Ic_Prioritys, therefore they alone cannot be cause of anomalies. Majority of anomalous incidents are lying above 'a' for the attribute Ir_OpenTime. According to it for all

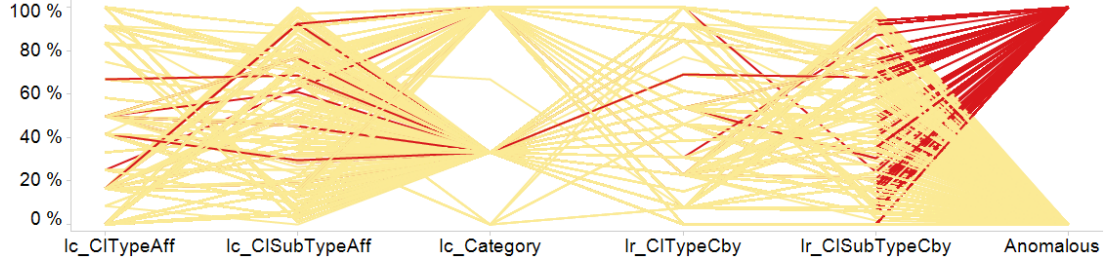


Figure 5.13: Parallel coordinate plot depicting the behaviour of incident CI type affected, incident CI subtype affected, incident category, interaction CI type affected and interaction CI subtype affected for anomalous and non-anomalous incidents from Markovian technique

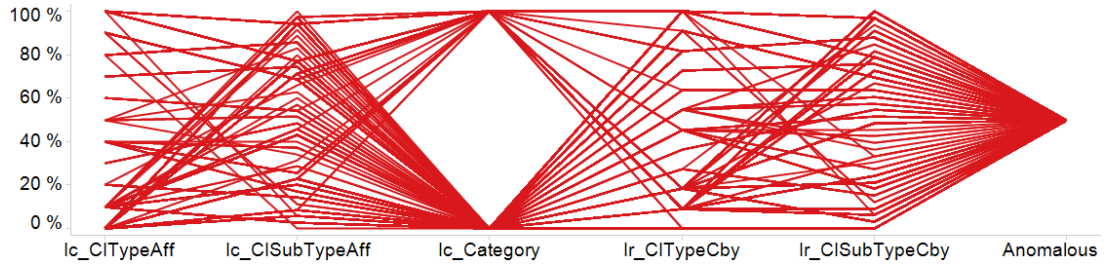


Figure 5.14: Parallel coordinate plot depicting the behaviour of incident CI type affected, incident CI subtype affected, incident category, interaction CI type affected and interaction CI subtype affected for anomalous incidents from Markovian technique

anomalous incidents, Ir_OpenTime is above 376305 (hrs) and majority of them have Ir_OpenTime above 381658 (hrs). Point 'b' shows that majority of anomalous incidents Ic_OpenTime above 381658 (hrs). Points 'c' and 'd' depicts that out of 4 Ic_Category : Incident, Complaint, Request for Information and Request for Change, anomalous incidents fall into only 2 categories: Incident and Request for Information.

Figure 5.12 shows parallel plot for the dataset obtained from table Interaction_Incident_Windows. We consider only records which have value of attribute Anomalous as 'YES'. Figure 5.12 shows that majority of anomalous incidents are lying above 'e' for the attribute Ir_OpenTime. According to it for all anomalous incidents, Ir_OpenTime is above 37082 (hrs) and majority of them have it above 381512 (hrs). Point 'f' shows that majority of anomalous incidents Ic_OpenTime above 381512 (hrs). Points 'g' and 'h' depicts that out of 4 Ic_Category's anomalous incidents fall into only 2 categories: Incident and Request for Information.

We choose attributes Ic_CITTypeAff, Ic_CISubTypeAff, Ic_Category, Ir_CITTypeAff

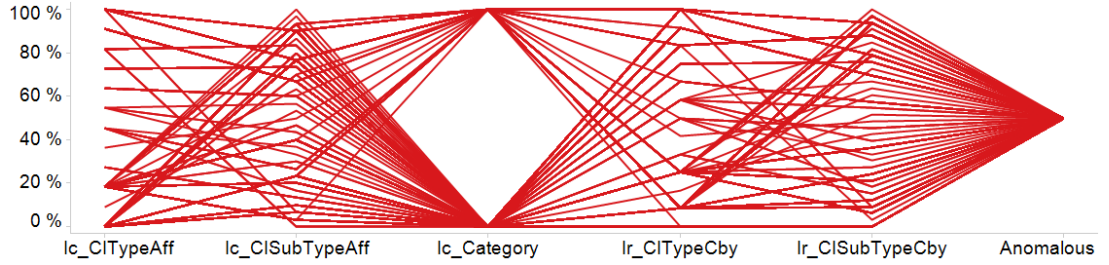


Figure 5.15: Parallel coordinate plot depicting the behaviour of incident CI type affected, incident CI subtype affected, incident category, interaction CI type affected and interaction CI subtype affected for anomalous incidents from Window based technique

and Ir_CISubTypeAff for Figure 5.13, 5.14 and 5.15. Figure 5.13 represents parallel plot for complete dataset with anomalous incidents obtained from Markovian technique. Figure 5.14 and 5.15 represents parallel plot only for anomalous incidents obtained from Markovian and Window based technique respectively. The selected attributes do not show any useful information regarding anomalous behavior of incidents. Anomalous Incidents are falling in all values for Ic_CITTypeAff, Ic_CISubTypeAff, Ir_CITTypeAff and Ir_CISubTypeAff. Therefore, they cannot be the cause of anomalies. It concludes that Affected Configuration Item's (CI) type and subtype do not influence cause of anomaly. The Affected Configuration Item is the CI where a disruption of ICT service is noticed. The attribute Ic_Category is showing the same behavior as in Figure 5.10, 5.11 and 5.12.

By comparing the parallel plots for Markovian and Window based technique, it is evident that anomalous incidents from both the techniques follow the same patterns.

5.6 Triangulation Study

In this Subsection, we present our approach on validating the uncovered root cause. In our experiments, we use a publicly available dataset and we do not have facts to validate the root-cause. The real source of the problem is confidential and not known to us or publicly available. We apply data triangulation technique¹ consisting of gathering evidences from multiple sources to validate the root cause [27]. Data triangulation is a well-known technique and we believe is well-suited for our study. We define two

¹[http://en.wikipedia.org/wiki/Triangulation_\(social_science\)](http://en.wikipedia.org/wiki/Triangulation_(social_science))

evaluators: outcome from parallel coordinate plots and output of decision trees on the dataset. Our objective is to investigate if the findings and indicators from the two different evaluators converge to the same conclusion. Decision tree results described in Section 5.4 show that root causes of anomalous incidents are open time of incidents alone or combination of open time of interaction, open time of incident, priority of incident and category of incident. Visualization using parallel plot depicts that cause of anomalies is not dependent on Affected Configuration Item's (CI) type and subtype which confirms with decision tree results. The parallel plot results also show that root cause of anomalies is dependent on open time of incident, open time of interaction and category of incident (refer to Section 5.5). The experimental results from the decision tree are in agreement with the parallel plot results, thereby validating our approach.

6

Limitations and Future Work

In our work, we have used a real world data from Rabobank Group. The size of dataset is small. The real source of the problem is confidential and not known to us or publicly available. Also anomalies are generally infrequent in nature. Therefore, in future we plan to validate our approach on bigger and other publicly available datasets¹². Also, in our work, we test our approach on IT Incident management domain. Future work includes applying our proposed approach for root cause analysis of anomalies in other domains.

We detect anomalies in business process logs which are sequential in nature. We use Markovian and Window Based technique for anomaly detection. In future, techniques like Kernal based and Hidden Markov Model based technique can also be used to detect anomalies. By using different techniques we can increase the confidence in results.

We use decision tree classifier to find root causes of anomalies. Future work includes testing methods like Naive Bayes, Support Vector Machine (SVM), Neural Networks, etc to find root causes.

We apply data triangulation technique consisting of gathering evidences from multiple sources to validate the root cause. All our analysis are based on two evidences of triangulation study that is from decision tree and parallel plot. We have used the publicly available dataset and we do not have facts to validate root-cause. In our future work, we would try to improve triangulation study by conducting experiments and gathering evidences with more classification and visualizations techniques.

¹<http://www.win.tue.nl/bpi/2013/challenge>

²<http://www.win.tue.nl/bpi/2012/challenge>

Conclusion

We present a novel approach for identification of anomalous traces and executions from event-logs generated by Process Aware Information Systems (PAIS) and a new technique for Root Cause Analysis (RCA) of anomalous traces. The key components of the proposed framework are: anomaly detection from sequential dataset using Window-based and Markovian-based technique, extraction of rules and characteristics of anomalous traces using decision-tree classifiers and application of parallel co-ordinate plots to visualize distinctions between anomalous and normal traces. We conduct a series of experiments on real-world dataset and conduct a triangulation study to demonstrate that the proposed approach is effective. Experimental results reveal agreement in output from Window-based and Markovian technique increasing the confidence in the classification result. We observe that data pre-processing and transformation is needed and impacts the outcome of parallel coordinate plot and decision tree classifier.

References

- [1] MARLON DUMAS, WIL M VAN DER AALST, AND ARTHUR H TER HOFSTEDE. *Process-aware information systems: bridging people and software through process technology*. John Wiley & Sons, 2005. 1
- [2] WMP VAN DER AALST. **Process mining: discovery, conformance and enhancement of business processes**. 2011. 1
- [3] WIL VAN DER AALST, TON WEIJTERS, AND LAURA MARUSTER. **Workflow mining: Discovering process models from event logs**. *Knowledge and Data Engineering, IEEE Transactions on*, **16**(9):1128–1142, 2004. 1
- [4] FÁBIO BEZERRA AND JACQUES WAINER. **Anomaly detection algorithms in logs of process aware systems**. In *Proceedings of the 2008 ACM symposium on Applied computing*, pages 951–952. ACM, 2008. 1, 4
- [5] V. CHANDOLA, A. BANERJEE, AND V. KUMAR. **Anomaly Detection for Discrete Sequences: A Survey**. *Knowledge and Data Engineering, IEEE Transactions on*, **24**(5):823–839, May 2012. 1, 7, 14, 15, 17, 19
- [6] GUILLERMO CALDERÓN-RUIZ AND MARCOS SEPÚLVEDA. **Automatic discovery of failures in business processes using Process Mining techniques**. 3
- [7] PETER TG HORNIX. **Performance analysis of business processes through process mining**. *Master’s Thesis, Eindhoven University of Technology*, 2007. 3
- [8] MITRA HERAVIZADEH, JAN MENDLING, AND MICHAEL ROSEMAN. **Root cause analysis in business processes**. 2008. 3

-
- [9] SURIADI SURIADI, CHUN OUYANG, WIL MP VAN DER AALST, AND ARTHUR HMTER HOFSTEDE. **Root cause analysis with enriched process logs**. In *Business Process Management Workshops*, pages 174–186. Springer, 2013. 4
- [10] ANDREAS ROGGE-SOLTI AND GJERGJI KASNECI. **Temporal Anomaly Detection in Business Processes**. In *Business Process Management*, pages 234–249. Springer, 2014. 4
- [11] EVGENIY VASILYEV, DIOGO R FERREIRA, AND JUNICHI IJIMA. **Using inductive reasoning to find the cause of process delays**. In *Business Informatics (CBI), 2013 IEEE 15th Conference on*, pages 242–249. IEEE, 2013. 4
- [12] STEPHEN MUGGLETON AND LUC DE RAEDT. **Inductive logic programming: Theory and methods**. *The Journal of Logic Programming*, **19**:629–679, 1994. 4
- [13] HENDRIK BLOCKEEL AND LUC DE RAEDT. **Top-down induction of first-order logical decision trees**. *Artificial intelligence*, **101**(1):285–297, 1998. 4
- [14] JACQUES WAINER, KWANGHOON KIM, AND CLARENCE A ELLIS. **A workflow mining method through model rewriting**. In *Groupware: Design, Implementation, and Use*, pages 184–191. Springer, 2005. 4
- [15] FÁBIO BEZERRA AND JACQUES WAINER. **Fraud detection in process aware systems**. *International Journal of Business Process Integration and Management*, **5**(2):121–129, 2011. 4, 5
- [16] FÁBIO BEZERRA, JACQUES WAINER, AND WIL MP VAN DER AALST. **Anomaly detection using process mining**. In *Enterprise, Business-Process and Information Systems Modeling*, pages 149–161. Springer, 2009. 4
- [17] FÁBIO BEZERRA AND JACQUES WAINER. **Algorithms for anomaly detection of traces in logs of process aware information systems**. *Information Systems*, **38**(1):33–44, 2013. 4, 5
- [18] FÁBIO BEZERRA AND JACQUES WAINER. **A Dynamic Threshold Algorithm for Anomaly Detection in Logs of Process Aware Systems**. *Journal of Information and Data Management*, **3**(3):316, 2012. 5

REFERENCES

- [19] S. FORREST, S.A. HOFMEYR, A. SOMAYAJI, AND T.A. LONGSTAFF. **A sense of self for Unix processes.** In *Security and Privacy, 1996. Proceedings., 1996 IEEE Symposium on*, pages 120–128, May 1996. 15
- [20] STEVEN A HOFMEYR, STEPHANIE FORREST, AND ANIL SOMAYAJI. **Intrusion detection using sequences of system calls.** *Journal of computer security*, **6**(3):151–180, 1998. 17
- [21] DANA RON, YORAM SINGER, AND NAFTALI TISHBY. **The power of amnesia: Learning probabilistic automata with variable memory length.** *Machine learning*, **25**(2-3):117–149, 1996. 17
- [22] BLAZ ZUPAN AND JANEZ DEMSAR. **Open-source tools for data mining.** *Clinics in laboratory medicine*, **28**(1):37–54, 2008. 21
- [23] IAN H WITTEN, EIBE FRANK, LEONARD E TRIGG, MARK A HALL, GEOFFREY HOLMES, AND SALLY JO CUNNINGHAM. **Weka: Practical machine learning tools and techniques with Java implementations.** 1999. 21
- [24] JOHN ROSS QUINLAN. *C4. 5: programs for machine learning*, **1**. Morgan kaufmann, 1993. 23
- [25] COHEN WILLIAM ET AL. **Fast effective rule induction.** In *Twelfth International Conference on Machine Learning*, pages 115–123, 1995. 23
- [26] YONGHENG ZHAO AND YANXIA ZHANG. **Comparison of decision tree methods for finding active objects.** *Advances in Space Research*, **41**(12):1955–1959, 2008. 23
- [27] JO MORAN-ELLIS, VICTORIA D ALEXANDER, ANN CRONIN, MARY DICKINSON, JANE FIELDING, JUDITH SLENEY, AND HILARY THOMAS. **Triangulation and integration: processes, claims and implications.** *Qualitative research*, **6**(1):45–59, 2006. 29