



HAL
open science

A Framework for Mining RFID Data From Schedule-Based Systems

Gurdal Ertek, Xu Chi, Allan N Zhang

► **To cite this version:**

Gurdal Ertek, Xu Chi, Allan N Zhang. A Framework for Mining RFID Data From Schedule-Based Systems. *IEEE Transactions on Systems, Man, and Cybernetics: Systems*, 2017, 47 (11), pp.2967-2984. 10.1109/TSMC.2016.2557762 . hal-01744328

HAL Id: hal-01744328

<https://hal.science/hal-01744328>

Submitted on 27 Mar 2018

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

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

Dr. Gürdal Ertek's Publications



A Framework for Mining RFID Data From Schedule-Based Systems

Gurdal Ertek, Xu Chi, Allan N. Zhang

Please cite this paper as follows:

Ertek, G., Chi, X., & Zhang, A. N. (2017). **A Framework for Mining RFID Data From Schedule-Based Systems**. IEEE Transactions on Systems, Man, and Cybernetics: Systems, 47(11), 2967-2984. DOI: 10.1109/TSMC.2016.2557762

Note: This document a draft version of this paper. Please cite this paper as above. You can download this draft version from the following website:

<http://ertekprojects.com/gurdal-ertek-publications/>

The published paper can be accessed from the following url:

<http://ieeexplore.ieee.org/abstract/document/7473879/>

A Framework for Mining RFID Data From Schedule-Based Systems

Gürdal Ertek, Xu Chi, *Member, IEEE*, and Allan N. Zhang *Member, IEEE*

Abstract—A *schedule-based system* is a system that operates on or contains within a schedule of events and breaks at particular time intervals. Given RFID data from a schedule-based system, what set of actions and computations, and what type of data mining methods can be applied so that one can obtain actionable insights regarding the system and domain? The research goal of this paper is to answer this posed research question through the development of a framework that systematically produces actionable insights for a given schedule-based system. We show that through integrating appropriate data analysis methodologies as a unified framework, one can obtain many insights from even a very simple RFID dataset, which contains only very few fields. The developed framework is general, and is applicable to any schedule-based system, as long as it operates under a few basic assumptions. The types of insights are also general, and are formulated in the most abstract possible way. The applicability of the developed framework is illustrated through a case study, where real world data from a schedule-based system is analyzed using the introduced framework. Insights obtained include the profiling of entities and events, the interactions between entity and events, and the relations between events.

Index Terms—Data mining, Decision support systems, Information systems.

I. INTRODUCTION

THE topic of this paper is the mining of data collected through RFID from schedule-based systems. A *schedule-based system* is a system that operates on (or contains within) a schedule of events and breaks at particular time intervals [1], [2]. Figure 1 illustrates a schedule-based system, which is characterized by a set of *entities* (or resources) \mathcal{I} entering and exiting a particular set of locations that have *events* \mathcal{J}^1 taking place in them according to a schedule. An entity is a distinct, independent, or self-contained being. An event is something that occurs in a certain place/location during a particular interval of time. The events may take place successively or may be separated by *breaks*, in other words, *time intervals* of no events. The set of breaks is denoted by \mathcal{J}^0 . Events and breaks constitute the set of time intervals \mathcal{J} . The schedule-based systems that we are particularly interested in are systems where the entry and exits of entities to location(s) are recorded through a data collection system, typically barcode, RFID, GPS (Global Positioning System), or sensors. Since RFID systems are gaining increasing importance in industry, we have illustrated a schedule-based system with RFID.

G. Ertek is with Rochester Institute of Technology - Dubai, Dubai Silicon Oasis, Dubai, UAE, e-mail: gurdalertek@gmail.com.

X. Chi and A.N. Zhang are with Singapore Institute of Manufacturing Technology, 71 Nanyang Dr, 638075, Singapore, e-mail: cxu@simtech.a-star.edu.sg, nzhang@simtech.a-star.edu.sg.

Manuscript received June 19, 2014; revised December 9, 2014.

Schedule-based systems are extensively encountered in a variety of domains, ranging from manufacturing to social event management. However, the basic elements of the system are the same. The basic elements are shown in bold in Figure 1. Table I lists some of the domains where schedule-based systems are present, and maps the key elements of a schedule-based system to domain-specific terminology.

An RFID (Radio Frequency Identification) system consists of tags (a.k.a. transponders) and readers (a.k.a. interrogators), typically also linked to an information system [3], [4]. In passive RFID, the information on the chip of the tag is read by the reader through radio waves, and the tag cannot transmit radio waves by itself. In active RFID, the tag has its own internal power source and the capability of actively transmitting information to the reader. Passive tags have the advantage of being significantly cheaper, whereas active tags possess larger memory capacity and can be used in more sophisticated scenarios. [3] provides an extensive review of RFID technology and its application in various industries, including logistics, retailing, travel and tourism, library science, food services and health care. A recent study reveals that only 3 percent of the companies in Europe have adopted RFID technology [3]. Thus, only a small percentage of companies have adopted RFID technology in their operations so far. However, the commitment of leading institutions (such as the US Department of Defense) and companies (such as Walmart, JC Penney and PG) is expected to eventually spread the use of RFID, just as the barcode technology has gained acceptance over time. [3], [5], and [6] provide a detailed discussion of RFID application domains, as well as a detailed literature review of RFID. [7] provides a highly useful list of potential benefits of RFID systems on operations management activities, in a multitude of domains. These benefits include preventing theft and shrinkage, identifying causes of spoilage, and evaluating employees.

RFID systems are used to basically produce data that can be mined through data mining methods for knowledge discovery and obtaining actionable insights. Data mining is the growing field of computer science where the goal is to uncover hidden information in -typically large and complex- piles of data [8]. There exist a multitude of data mining methods that can be applied depending on the size and structure of the data at hand. Data mining can thus be considered as a field which encompasses a collection of interrelated and interacting tools, including clustering, classification, association mining, network analysis, data visualization, as well as others. A significant challenge then is the selection of the appropriate set of methodologies and the way they are applied in analyzing a

particular dataset.

The research question to be answered in this paper is the following:

“Given RFID data from a schedule-based system in any domain (such as social event management, manufacturing, healthcare, etc.) what set of actions (including the data cleaning steps) and computations, and what type of data analysis and data mining methods can be applied, so that one can obtain actionable insights regarding the system and the domain?”

The research goal to answer the above research question is the development of a framework, that takes RFID data and basic event schedule data and information, and produces actionable insights regarding the system and entities within the system. Our first main motivation was to show that, through appropriate data analysis methodologies, one can obtain many insights from even a very simple RFID dataset, which contains only very few fields. Our second main motivation was that such a framework would be applicable in a wide range of domains. Our third motivation was observing from our survey of the literature that there is a significant gap regarding this type of research.

The contributions of our study are multifold: First, we introduce an analysis framework, including its mathematical representation, for mining RFID data coming from a schedule-based system. The framework developed is general, and is applicable to any schedule-based system that operates as described. While the framework is developed assuming a single location, it can also be extended to the case of multiple locations by introducing a set of locations \mathcal{L} and a new dimension in the relevant sets and parameters. Second, we enumerate the different types of insights that can be obtained through the introduced framework. These insights are also general, and are formulated in the most abstract way possible. Third, we develop and present the corresponding algorithms that are needed in the analysis framework. The framework depends on these algorithms to do the required data processing, database augmentation, and other computations. Finally, we demonstrate the applicability of the developed framework through a case study, where real world data from a schedule-based system is analyzed using the introduced framework. The case study illustrates how the framework can be applied in the real world for a given domain.

The novelty of the research is the introduction of a data mining framework for the first time for this type of a system. The existing research in schedule-based systems mainly focuses on obtaining good, and if possible optimal, schedules, or event processing. However, the interaction of the entities in the system, given the obtained schedule, has not been analyzed in depth in earlier research. The importance of the research lies in its general applicability in a wide range of domains. Table I lists some of the application areas of the developed framework, with a mapping to the domain-specific terminology. Thus the developed framework is applicable in its current form in all the listed domains, because the fundamental aspects of the model are the same across domains.

The remainder of the paper is organized as follows: Section II provides a brief review of some relevant literature as the

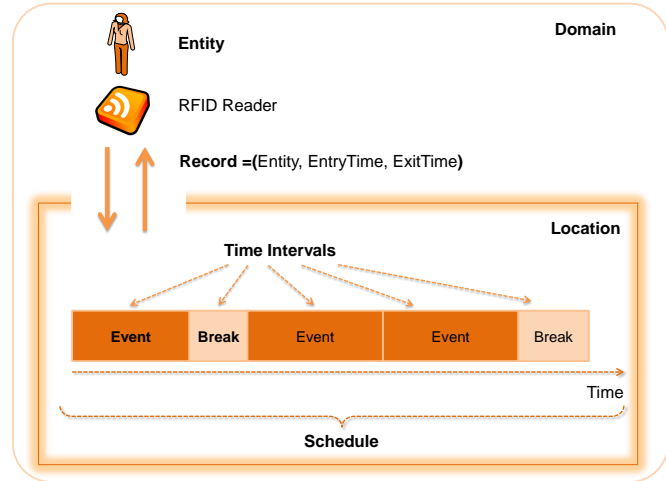


Fig. 1. A schedule-based system where entities entering and exiting the system are tracked with RFID

Domain	Entity	Location	Event	Break
Manufacturing	Worker, Cart	Workshop, Cell	Production, Setup	No production, Break
Warehousing	Worker, Pallet, Forklift	Warehouse, Zone	Order Pick Wave	No picking, Break
Transportation	Pallet	Truck, Warehouse	Transportation, (Un)Loading	Stalling, No operation
Social Event	Attendee, Organizer	Room, Hall	Session	Break
Healthcare	Doctor, Nurse, Patient, Equipment	Surgery Room, Hospital rooms	Surgery, Test Appointment	Break
Education	Instructor, Student	Classroom, Laboratory	Course hours, Labs, Recitations	Break
Tourism	Tourist, Visitor	Museum, Art Gallery	Guided tour, Activity	No tour

TABLE I
APPLICATION AREAS OF THE DEVELOPED METHODOLOGY, WITH A MAPPING OF THE VARIOUS ASPECTS OF THE MODEL.

background. Section III discusses the framework developed and proposed. Section IV is devoted to the results and analysis of the case study, where new insights are obtained. Finally, Section VI presents some conclusive remarks.

II. LITERATURE

A. Schedule-based Systems

The primary line of existing research regarding schedule-based systems involves the derivation of good, and if possible optimal, schedules. The primary modeling approach for this line of research is optimization, and typically mixed-integer programming. [9] is the classic reference for scheduling theory, and [10] contains a detailed discussion of practice and application of scheduling, in addition to theory and algorithms. The scheduling research focuses on whether problems are polynomially solvable and optimal under certain conditions [11]. Typical contribution in such research also includes optimization or approximation algorithms and analysis of worst case error bound. Scheduling can be at any resolution, ranging from single-machine machine scheduling [11] to the scheduling of supply chains [12]. One line of scheduling research develops or applies machine learning and data mining methods and algorithms for generating the schedules [13]. Some of these studies also analyze generated schedules using data mining techniques for coming up with new schedules

[14]–[16]. However, while very extensive research exists on scheduling, the interaction of the entities in the system, given the obtained schedule, has not been analyzed from a data mining perspective in earlier research. In our research, we provide the possible practical benefits of such a perspective in Section V. One final stream of research regarding schedule-based systems is regarding the processing of the events data [17].

B. Mining RFID Data

There exists a large body of literature on the mining of RFID data. However, an extensive survey performed during our study revealed that none of the existing research studies have developed a comprehensive framework for mining RFID data coming from a schedule-based system. One approach could be modifying Knowledge Discovery and Data Mining (KDDM) process models [18], [19] for this particular domain.

The most time consuming step in data mining is typically data cleaning. [20] develops a framework for RFID data cleaning. [21] presents a data cleaning methodology for indoor RFID data, eliminating temporal redundancy and spatial ambiguity, by building a distance-aware graph. The authors test and illustrate the methodology with real data from the baggage handling system of an airport.

The success of a data mining process is highly dependent on the underlying data structure. To this end, [6] develops a data mining infrastructure that allows the efficient data mining of RFID data. Specifically, the authors introduce two new data models, namely path cube and workflow cube. They explain and illustrate their approach using examples and data from supply chain management.

Based on our literature review, the domains where one can find the mining of RFID data are supply chain management and logistics, as well as retail.

[22] presents a data processing and mining framework for logistics using RFID data. [23] performs a rule-based analysis and GIS-based visualization of RFID data for managing items in a supply chain. For example, consistency of velocity and waiting time has to be ensured for an item throughout the supply chain, and any anomalies have to be detected. In a similar study, [24] applies 3-dimensional visualization for tracking and understanding object movements through time, again enabling the discovery of irregularities.

The following three studies are examples of data mining for retail RFID data: [25] develops a framework for the analysis of residence time in shopping, based on the mining of RFID data. [7] and [26] use RFID data for targeted advertising inside a retail store. [27] uses RFID data for predicting retail store sales.

Studies on the mining of RFID data for other domains include the following: [28] presents a framework for quality assurance, as well as two industry applications. [29] mines RFID data through the integration of fuzzy logic for resource allocation in garment manufacturing, and illustrates the applicability of this approach at a company. [30] presents a framework that uses RFID data for intelligent traffic management. [31] performs sequential pattern mining of RFID data for

generating tourist path suggestions. [27] recommends routes for theme park visitors using real time RFID information and historical tourist behavior data. [32] presents a knowledge-based system framework for healthcare using RFID data. [33] mines RFID data for smart home prediction.

A multitude of studies investigate the outlier detection problem with RFID data. [34] carries out behavior modeling using RFID data, using clustering to detect abnormal events. [35] also performs behavior modeling using RFID data, detecting abnormal events in elderly care. [36] uses RFID data for behavior identification and anomaly detection. [37] mines frequent trajectory patterns and detects abnormal trajectories. [38] presents a data mining framework that detects outlier observations in RFID data. [39] analyzes the dynamics of person-to-person interaction networks using RFID data.

Other related papers do not necessarily use data collected through RFID, but illustrate methods and case studies that can be adopted to the analysis of RFID based data. For example, [40] presents a very detailed analysis of data on location-based social networks. Some of the research questions investigated in [40] include how social connection is affected by geographical distance, how users can be clustered based on their activities, how user mobility is influenced by various factors, and how home locations of users can be predicted. [41] mines matching behavioral patterns based on joining various kinds of entity characteristics in mobile communication. One final related line of research builds social recommender systems with various benefits, such as supporting the creation of new social relations [42].

C. Mining RFID Data from Social Events

RFID technology has a great potential for facilitating and enhancing the management of social events, where humans interact with each other over time and across different locations. The case study in our paper presents the application of RFID in the context of a social event, specifically a scientific conference. [43]–[46] provide information system architectures for collecting data in a conference through RFID. [46] also describes how this data can be used in real time for informing conference attendees and illustrates, through a detailed scenario, how the system operates. [47] describes how UML (Unified Modeling Language) can be extended to model Web 2.0-based context-aware applications. The UML profile explained in [47] can be used in developing the browser-accessed online services and mobile applications that can be deployed for conference management.

RFID systems, when used in social events, generate time-stamped location data for each of the attendees/participants of the event. This data, when combined with other data regarding the attributes of the attendees, locations and the event schedule, can generate significant insights regarding the attendees, the structure and the nature of the social network, and the event. Furthermore, the methods employed for mining social network data [48], [49] can be fused to obtain hybrid data analysis frameworks. These insights and the information systems designed around them can be used to improve the social event in better serving its intended goals [45]. Improved

conference management information systems and managerial practices can enable the attendees find sessions and other people that they would be interested in, minimize schedule conflicts, increase participation in the sessions, and improve the overall quality of the event.

[50] integrates RFID data with online data from social networks (e.g. Facebook, Twitter) and offline data from earlier conferences, and develops an ubiquitous conference management system. The system generates context-aware recommendations to conference attendees, significantly increasing attendees' satisfaction with the event.

[43], [44], and [51] are the most related studies in the literature to our case study, because these papers carry out posterior visualization and analysis of RFID enriched event data, and furthermore give examples of insight-generating questions whose answers can be obtained through querying the data.

[43] develops an infrastructure and a scalable information system for tracking and analyzing human face-to-face (f2f) contact networks, such as people in a scientific conference. The authors employ RFID technology and data reporting and analysis methods for enhancing social interactions between event attendees and industry exhibitors.

[44] develops an RFID based system for connecting conference attendees based on their locations, the sessions that they have attended, and the attendees they have interacted with. Posterior analysis of attendee behavior suggested that earlier physical encounter during the conference (proximity), as well as commonality of attributes (homophily) were the most important factors affecting the selection of new contacts.

[51] also presents an information system for conference management and detailed analysis of the obtained f2f data, as in [43] and [44]. The main contribution of [51] is the comprehensive evaluation of the behavioral patterns in a conference setting, developing analysis techniques for revealing roles of the attendees and attendee communities. Explicit and organizing roles are discovered through the analysis of classic centrality measures used in graph theory, such as degree, strength, betweenness, closeness, and eigenvalue centrality.

While [43], [44], and [51] bring fresh perspectives to the mining of RFID data, our work has several additional aspects in comparison to these studies: First, we develop a complete framework that exhaustively explores and exhibits all the possible types of insights, rather than a set of selected few insights. Second, our framework requires a very basic data, with very few attributes, collected by almost every RFID system by default. Third, our framework is described not only conceptually, but also through rigorous mathematical formalism. The algorithms used for data processing are also included in the work. Fourth, rather than discussing a single domain, we generalize the analysis to schedule-based systems, which can include a very rich collection of application domains. Fifth, we discuss the practical implications of our research for not only a single domain (ex: social event management), but for a multitude of domains.

III. FRAMEWORK

In this section, we describe the framework that we introduce for mining RFID data from schedule-based systems. First we outline the research steps followed in the study. Then we list our assumptions regarding the analyzed system. Third, we introduce the mathematical notation and the database structures in the various stages of the analysis framework. Fourth, we describe the computational algorithms for augmenting the RFID data obtained from a schedule-based system. Finally, we present the novel analysis framework that we have developed, and list the types of insights that can be obtained through this framework.

A. Research Steps

Our study consists of the steps listed below, and resulted in the framework and case study presented in this paper. We thus suggest the application of similar steps in analyzing the RFID data coming from a system with particular characteristics.

- 1) Understanding of the data mining research goal, as well as the research question and the domain.
- 2) Development of a mathematical notation (example: sets, parameters,...)
- 3) Description of the RFID data and the domain-related data in terms of the developed mathematical notation (example: entities, entity entrance times to events)
- 4) Identification of the metrics to be computed (example: whether an attendee has attended a particular session or not, as well as the time s/he spent in each session), and the database structures needed.
- 5) Development of formulas for obtaining the desired performance metrics and insights.
- 6) Identification some of the possible types of data analysis that can be implemented, as well as some of the possible types of insights that can be obtained through each type of data analysis.
- 7) Survey of the literature for related studies and recording the types of analysis they present, which can be adopted.
- 8) Execution of the data analysis process, and the discovery of various types of insights.
- 9) Elicitation of the obtained results and insights, and the subsequent filtering of the most essential and actionable insights among those obtained. The importance and actionability of insights were decided upon through discussion sessions with conference organizers from academia.
- 10) Integration of the executed data mining processes in a single unified framework, and proposing it as a general methodology for the analysis of RFID data from schedule-based systems, that can be applied to systems other than schedule-based ones.

B. Assumptions

Our assumptions regarding the RFID data collection are as follows:

- 1) The gateway where the RFID reader is located is an in-out-gateway [52].

- 2) RFID tags are read throughout the event schedule, not missing any of the events, nor people passing through the doors.
- 3) All passes (entries and exits) made with an RFID tag are read, with the RFID receiver not missing any passes.
- 4) RFID readings and the final data are accurate.
- 5) Every entity wears RFID during passes, except when the RFID tag is left in the location, never to be worn again.
- 6) All events happen in one location.

These assumptions (except the last) are required so that the data is accurate and complete. The last assumption is assumed so that the concepts and the developed framework can be easily demonstrated.

C. Mathematical Notation

We now introduce the mathematical notation that will be used throughout the description of the framework. While the indices are always provided in the notation, for convenience, sometimes the indices are dropped (for example, u). In that case, the symbol refers to the symbol with the default indices that were specified when the notation was initially introduced (for example, u refers to u_{ir} , because that is how it is defined initially). The database structures and algorithms will also be introduced in this subsection.

Sets

\mathcal{R} : set of unique record IDs ; $r : 1 \dots R$

\mathcal{I} : set of entities ; $i : 1 \dots I$

\mathcal{J} : set of time intervals ; $j : 0, 1 \dots J$ (The time intervals correspond to actual events and the breaks between these events); $\mathcal{J} = \mathcal{J}^0 \cup \mathcal{J}^1$.

\mathcal{J}^0 : set of breaks

\mathcal{J}^1 : set of events

Given Data

u_{ir} : entry time of entity i in record r

U_{ir} : exit time of entity i in record r

\mathcal{D}^0 : the database of RFID logs; $\mathcal{D}^0 = \{d^0 : \langle r, i, u_{ir}, U_{ir} \rangle\}$

Event Schedule Data

s_j : start time of time interval j

f_j : finish time of time interval j

d_j : duration of time interval j ; $d_j = f_j - s_j$

\mathcal{D}' : the database of time intervals;

$$\mathcal{D}' = \{d' : \langle j, intervalType(j), s_j, f_j, d_j \rangle\}$$

where $intervalType(j)$ is a lookup function (defined next) that returns whether the time interval corresponds to an event or a break.

Lookup Functions

$$intervalType(j) = \begin{cases} break, & \text{if } j \in \mathcal{J}^0 \\ event, & \text{if } j \in \mathcal{J}^1 \\ null, o/w \end{cases}$$

$$intervalOf(t) = \{j \in \mathcal{J} : s_j \leq t \leq f_j\}$$

Intermediary Data

$u = u_{ir}$: entry time of an entity in a record

$U = U_{ir}$: exit time of an entity in a record

$e = e_{irj}$: entry time of an entity to an event in a record

$x = x_{irj}$: exit time of an entity from an event in a record

$T = T_{irj} = x - e$: time spent by entity at an event (in a single record)

$p = p_{irj}$: start time of an entity present at the location for an event in a record (The entity may wait for the event.)

$q = q_{irj}$: end time of an entity present at the location for an event in a record (The entity may be spending additional time at the location after the event is completed.)

Computed Metrics

\underline{p}_{ij} : earliest start time of an entity present at the location for an event; $\underline{p}_{ij} = \min_{r \in \mathcal{R}} p_{irj}$.

\bar{q}_{ij} : latest end time of an entity present at the location for an event; $\bar{q}_{ij} = \max_{r \in \mathcal{R}} q_{irj}$.

$earliness' = earliness_{irj}$: how early entity i entered event j in a given record r ; takes positive value if entity entered early, and takes negative value if entity entered late; $earliness_{irj} = s_j - p_{irj}$.

$lateness' = lateness_{irj}$: how late entity i exited from event j in a given record r ; takes positive value if entity exited late, and takes negative value if entity exited early; $lateness_{irj} = q_{irj} - f_j$.

$earliness = earliness_{ij}$: how early entity i entered event j ; takes positive value if entity entered early, and takes negative value if entity entered late; $earliness_{ij} = \max_{r \in \mathcal{R}} earliness_{irj}$.

$lateness = lateness_{ij}$: how late entity i exited event j ; takes positive value if entity exited late, and takes negative value if entity exited early; $lateness_{ij} = \max_{r \in \mathcal{R}} lateness_{irj}$.

$$entryStatus = \begin{cases} NoEntry, & \text{if } earliness = null \text{ and not an entry from previous event} \\ EarlyEntry, & \text{if } earliness \geq 0 \text{ and not an entry from previous event} \\ LateEntry, & \text{if } earliness < 0 \text{ and not an entry from previous event} \\ EntryFromPreviousEvent, & \text{if entry from previous event} \end{cases}$$

$$exitStatus = \begin{cases} NoExit, & \text{if } lateness = null \text{ and not an exit into next event} \\ EarlyExit, & \text{if } lateness \geq 0 \text{ and not an exit into next event} \\ LateExit, & \text{if } lateness < 0 \text{ and not an exit into next event} \\ ExitIntoNextEvent, & \text{if exit into next event} \end{cases}$$

$$Z_{ij} = \begin{cases} 1, & \text{if attendee } i \text{ attended event } j, i \in \mathcal{I}, j \in \mathcal{J} \\ 0, o/w \end{cases}$$

n_{ij} : number of times that entity i has entered and/or exited event j

T_{ij} : total time (stay duration) that entity i spent in event j ;

$$T_{ij} = \sum_{r \in \mathcal{R}} T_{irj}$$

Databases

The databases whose structures are given here are shown as cylinders in Figure 2. For example, cylinder with the label 0 refers to \mathcal{D}^0 and the cylinder with the label 1 refers to \mathcal{D}^1 .

The database structures for the Raw RFID Database and the Joined Database 1 are

$$\mathcal{D}^0 = \{d^0 : \langle r, i, u_{ir}, U_{ir} \rangle\}$$

$$\mathcal{D}^1 = \{d^1 : \langle r, i, u_{ir}, U_{ir}, j_1, j_2, \varepsilon_1, \varepsilon_2, s_{j_1}, s_{j_2}, f_{j_1}, f_{j_2} \rangle\}$$

The Augmented Database is

$$\mathcal{D}^2 = \{d^2 : \langle i, j, r, u, U, e, x, t, p, q, \text{earliness}', \text{lateness}', \text{entryStatus}, \text{exitStatus} \rangle, j \in \mathcal{J}^1\}$$

Entity-Event Profile database and the databases derived from that database are

$$\mathcal{D}^3 = \{d^3 : \langle i, j, \text{earliness} \rangle, j \in \mathcal{J}^1\}$$

$$\mathcal{D}^4 = \{d^4 : \langle i, j, \Gamma_1(\text{earliness}) \rangle, j \in \mathcal{J}^1\}$$

where

$$\Gamma_1(\text{earliness}) = \begin{cases} \text{NoEntry}, & \text{if } \text{earliness} = \text{null} \\ \text{EarlyEntry}, & \text{if } \text{earliness} \geq 0 \\ \text{LateEntry}, & \text{if } \text{earliness} < 0 \end{cases}$$

$$\mathcal{D}^5 = \{d^5 : \langle i, j, \Gamma_2(\text{earliness}) \rangle, j \in \mathcal{J}^1\}$$

where

$$\Gamma_2(\text{earliness}) = \begin{cases} \text{NoEntry}, & \text{if } \text{earliness} = \text{null} \\ \text{Entry}, & \text{o/w} \end{cases}$$

Entity Profile database is

$$\mathcal{D}^6 = \{d^6 : \langle i, \text{avg}(t), \text{avg}_j(\text{earliness}), \text{min}_j(\text{earliness}), \text{max}_j(\text{earliness}), \text{stdev}_j(\text{earliness}), \text{avg}_j(\text{lateness}), \text{min}_j(\text{lateness}), \text{max}_j(\text{lateness}), \text{stdev}_j(\text{lateness}), \text{count}_j(i) \rangle, j \in \mathcal{J}^1\}$$

where the computations for \mathcal{D}^6 are done using \mathcal{D}^2 .

The remaining databases are

$$\mathcal{D}^7 = \{d^7 : \langle i_1, i_2 \rangle, i_1, i_2 \in \mathcal{I}\}$$

$$\mathcal{D}^8 = \{d^8 : \langle i_1, i_2, \text{count}(i_1, i_2) \rangle, i_1, i_2 \in \mathcal{I}\}$$

$$\mathcal{D}^9 = \{d^6 \cup \text{metrics}(i), i \in \mathcal{I}', d^6 \in \mathcal{D}^6\}$$

where \mathcal{I}' is the set of entities in \mathcal{D}^8 which have a support count greater than the minimum support count threshold, and $\text{metrics}(i)$ is a function that returns the array of computed graph metrics for an item i .

D. Computational Algorithms for Data Augmentation

The computational algorithms for augmenting the RFID data are given in Appendix A. The first of these algorithms takes the raw RFID data \mathcal{D}^0 and the schedule data, and joins these two tables to form a new data table, namely \mathcal{D}^1 . The second algorithm is more complicated, and is focused only in what is happening with respect to events (rather than breaks). This second algorithm transforms the data which is in the form of entry/exit records into a database \mathcal{D}^2 that contains information only on entities and events. The augmented database \mathcal{D}^2 contains the entry and exit times of entities to events, as well as their earliness (positive value if early entry), lateness (positive value if late exit), as well as other data. \mathcal{D}^2 is critical, because it is used in later stages of the analysis framework to extract new databases and to obtain insights.

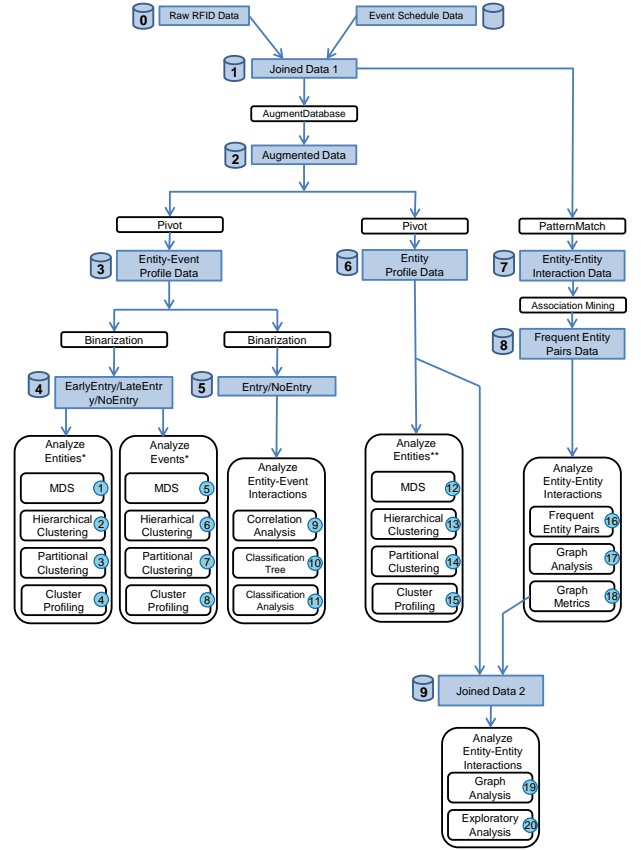


Fig. 2. The developed analysis framework for mining RFID data from a schedule-based system

E. Analysis Framework

The developed analysis framework is given as a flowchart in Figure 2. The framework starts with raw data coming from RFID system, as well as data regarding the schedule of events in the system. The data is then brought to a richness so that it can be analyzed to obtain insights. The analysis centers around three lines; shown with the numbers 3, 6, and 7 in the figure. The insights are obtained through analyzing entities, events, entity-event interactions, and entity-entity interactions.

Figure 2 shows that the analysis begins by joining the raw RFID data \mathcal{D}^0 (shown with the cylinder with the label 0) with the event schedule data to form \mathcal{D}^1 , and then augmenting \mathcal{D}^1 to generate \mathcal{D}^2 . Next, three basic types of data are obtained:

\mathcal{D}^3 , *Entity-Event Profile Data* is obtained through pivoting on \mathcal{D}^2 , and shows the earliness of each entity for each event. Some of the values are missing, indicating that the entity did not enter the system at all during a particular event.

\mathcal{D}^6 , *Entity Profile Data* shows the metric statistics for each entity as computed over all the sessions.

\mathcal{D}^7 , *Entity-Entity Interaction Data* lists the entity pairs that have entered or exited the system simultaneously. The data is obtained through running a pattern matching algorithm (Appendix B).

Having obtained these three basic types of data, further data transformations and/or algorithms are applied to obtain

Insight No	Question Answered	Example in
	Behavior Analytics 1 Based on the behavioral patterns of the entities with respect to specific events...	
1	Which entities are positioned close to each other?	Figure 3
2	Which groups of entities behave most similar?	
3	Which entities can be clustered into which clusters?	
4	What are the profiles of these entity clusters?	
	Event Analytics Based on the behavioral patterns of the entities with respect to specific events...	
5	Which events are similar to each other?	
6	Which groups of events are most similar?	Figure 4
7	Which events can be clustered into which clusters?	
8	What are the profiles of these event clusters?	
9	What is the correlation between different events?	Table III
10	Which earlier events affect a particular event, and how?	Figure 5
11	Can the entry of specific entities to an event be predicted?	Table IV
	Behavior Analytics 2 Based on the general behavioral patterns of the entities...	
12	Which entities are positioned close to each other?	Figure 6
13	Which groups of entities behave most similar?	Figure 7
14	Which entities can be clustered into which clusters?	
15	What are the profiles of these entity clusters?	Figure 8
	Relationship Network Analysis Based on the joint actions of the entities...	
16	Which entity pairs enter/exit many events together?	Table V
17	How are the entities related to each other?	Figure 9
18	Which entities are influencers and which are followers?	Table VI
19	How can the behavioral attributes be analyzed together with the relationship network?	Figure 10
20	How can the behavioral attributes be analyzed together with the graph metrics?	

TABLE II

INSIGHTS THAT CAN BE OBTAINED THROUGH THE INTRODUCED ANALYSIS FRAMEWORK

insights into the system. These insight types are numbered from 1 to 20 in Figure 2, inside the circles. These 20 insight types are then listed in Table II. In Table II, the insights that can be obtained using our proposed framework have been classified into categories of behavior analytics, event analytics, and relationship network analysis. The behavior analytics category has been further labeled as 1 or 2 based on the analytics quest and the data source.

Table II also lists (in its last column) the figure and/or tables which illustrate the insight type in the case study. For example, consider the line corresponding to Insight 2 in Table II. This insight aims at answering the question “Which entities are positioned close to each other?”. The same line in the table tells that an example of the analysis that leads to this type of insight is illustrated in Figure 3. Due to space limitations in the paper, we are able to provide examples for only some of the insights, hence the empty cells under the last column of the table.

IV. CASE STUDY

In this section we firstly describe the data used in the case study, and then illustrate the various insights that can be obtained through the presented framework. The insights that will be presented are listed in Table II, along with the figure and table numbers. The data mining processes applied are described in Appendix C, and the software tools used are described in Appendix D.

A. Data

The data used in the study belongs to the domain of social event management, and comes from a four-day medical conference. Each attendee of the conference was provided with a unique RFID tag and their entry and exit times to the single conference hall were recorded. The raw RFID data D^0 consists of 9624 rows (entry-exit combinations), and four columns,

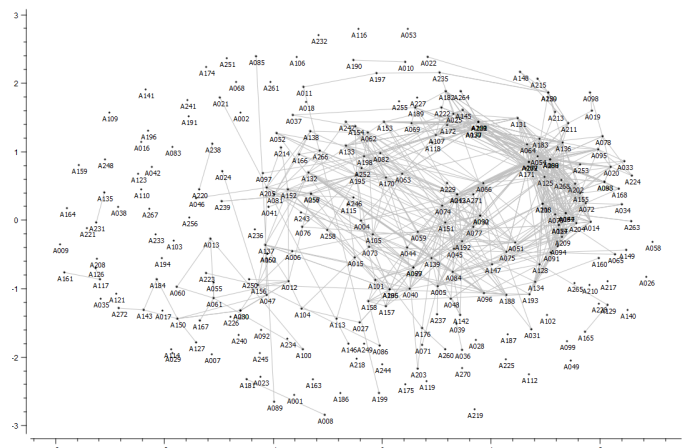


Fig. 3. MDS results for entity analysis based on entity-event profile data

where the columns are the record ID, attendee name (masked), entry date and time, and exit date and time. The total number of attendees (total number of entities in the schedule-based system) is 272. The schedule consists of 17 events and 11 breaks (including the time interval before the first event and the time interval after the last event) giving a total of 28 time intervals.

B. Behavior Analytics 1: Based on Entity-Event Data

The first illustration is for Insight 1, and is given in Figure 3. This insight answers the question “Which entities are positioned close to each other?”. The analysis here is based on temporal proximity [53] of entities. The data mining method used for this purpose is Multi-Dimensional Scaling (MDS) (as read from the box that corresponds to Insight 1 circle in Figure 2), which maps multi-dimensional data onto two dimensions, based on how close the data points are [54]. Figure 3 shows the mapping of attendees (entities) on a two-dimensional plane. The most significant associations are shown with lines between the points. Since Insight 1 is using the database D^4 (EarlyEntry/LateEntry/NoEntry data), the closeness of the points, as well as the links between them, are based on the Hamming distance in between. Hamming distance is a distance measure that computes the number of bits two strings are different from each other [55]. As an example, if two entities entered all events early, but differed only in their behavior with respect to one event (for example one entered early, and the other entered late into the last event), the Hamming distance between them would be 1. The Hamming distance measure has been selected, rather than other distance measures, because it is a very popular distance measure when the data points are binary vectors. One can observe a highly dense region in Figure 3, to the right of the figure, as well as a less dense region in the middle of the figure, and some sparse links. This shows that the entities in the dense cluster are very much close to each other, whereas there are other closely positioned entities among the remaining entities. Furthermore, given a particular entity, one can find the other entities closely positioned to this entity from the figure.

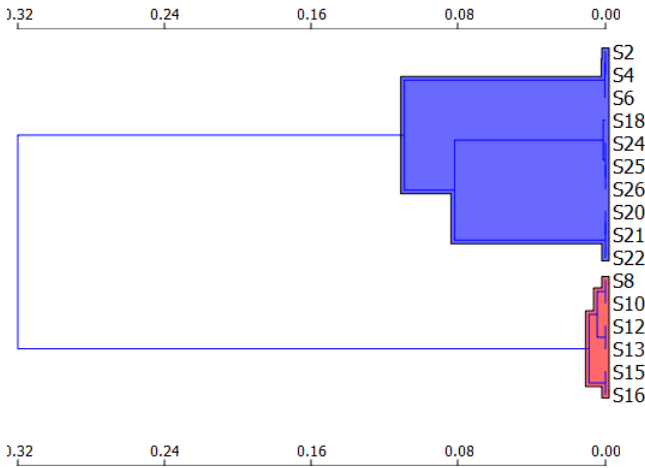


Fig. 4. Attribute (session) clustering results for Case 2

Session1	Session2	Are Successive?	Correlation
S25	S26	Yes	0.88
S20	S21	Yes	0.68
S2	S4	Yes	0.59
S25	S27		0.53
S26	S27	Yes	0.60
S21	S22	Yes	0.53
S24	S27		-0.50
S8	S10	Yes	-0.55
S15	S16	Yes	-0.66
S24	S26		-0.72
S24	S25	Yes	-0.77

TABLE III
HIGHEST AND LOWEST CORRELATIONS BETWEEN SESSIONS

C. Event Analytics Based on Entity-Event Data

The next illustration is for Insight 6, and is given in Figure 4. This insight answers the question “Which groups of events are most similar?”. The data mining method used for this purpose (as read from Figure 2) is hierarchical clustering, which hierarchically builds clusters from data, starting from individual points ([56], p44). An interesting point here is that the data points are not the entities, but rather the *events* (sessions in the case study). So the goal is to see which events are similar to each other, based on the entry-exit patterns of the entities. This analysis also uses database D^4 (EarlyEntry/LateEntry/NoEntry data), however, carries hierarchical clustering of *events* (rather than the entities), based on the Hamming distances between the events. The dendrogram in Figure 4 shows that events {S2, S4, S6} are similar to each other, based on the entity-event profile data. Similarly, {S18, S24, S25, S26} are very similar, since they form a cluster together. Other groups of similar events are {S20, S21, S22}, {S8, S10}, {S12, S13}, and {S15, S16}.

The next illustration is for Insight 9, and is given in Table III. This insight answers the question “What is the correlation between different events?”. The data mining / statistics method used for this purpose is correlation analysis, which computes the linear association between pairs of observations [57]. Again, the observations are events (rather than entities). Table

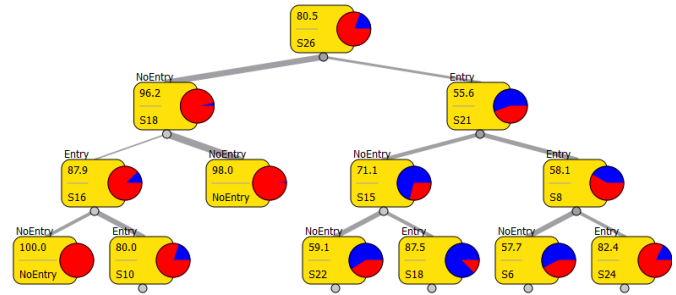


Fig. 5. The classification tree for the case where class attribute is entry into session S27

III shows the correlation values above 0.50 and below -0.50. High correlation between successive events indicates that the entities which entered the former of those successive events also mostly entered the latter, or those who did not enter the former did not enter the latter. This is the case for session pairs (S25, S26), (S20, S21), (S2, S4), (S26, S27), and (S21, S22). This high correlation can indicate that the former event encouraged entry to the latter, that the two events catered to the same set of entities, or both of these. Negative correlation between two successive events is also an important observation, and may be due to one or combination of several reasons: First, it may be that the former event was (not) successful, en(dis)couraging entry to the latter. Second, the two events may be catering to different set of entities. Third, there may be another reason, such as the latter event being the last event of the day, and entities exiting the system early. The successive event pairs with high negative correlation are (S8, S10), (S15, S16), and (S24, S25).

The following illustration is for Insight 10, and is given in Figure 5. This insight answers the question “Which earlier events affect a particular event, and how?”. The data mining method used for this purpose is classification tree analysis [58]. Classification trees summarize rule-based information about classification as trees. In classification tree models, each node is split (branched) according to a criterion. Then, a tree is constructed with a depth until all the rules are displayed on the graph under a stopping criterion. At each level, the attribute that creates the most increase compared with the previous level is observed. The algorithms for classification tree analysis are explained in [58]. In the implementation that we utilized in our analysis, selecting the attributes for the splits is based on information gain. In classification trees, identifying the nodes that differ noticeably from the root node are important, because the path that leads to those nodes tells us how significant changes are observed in the sub-sample compared with the complete data. By observing the shares of slices and comparing with the parent and root nodes, one can discover interesting classification rules and insights. Figure 5 shows the classification tree where the Entry/NoEntry into event S27 (last session in the case study) is the predicted attribute. The very first split, based on the value of S26 provides the most information. In the complete data, 80.5% of the entities did not participate in event S27 (light shaded slice). However, among

Classifier	CA	AUC
CN2 rules	0.8347	0.7925
kNN	0.8186	0.7763
Classification Tree	0.8327	0.7267
SVM	0.8274	0.8386
Naive Bayes	0.8243	0.8573
Neural Network	0.8315	0.8432

TABLE IV

THE CLASSIFICATION RESULTS FOR PREDICTING ATTENDEES TO SESSION S27 (LAST EVENT IN THE SCHEDULE)

those entities that did not enter S26 at all (NoEntry), this percentage is 96.2%. On the other hand, among the entities that did enter S26, the percentage of Entry into S27 is higher (55.6%). So, approximately half (55.6%) of those entities that entered S26 entered S27, whereas almost all (96.2%) of the entities that skipped S26 also skipped S27. This connectedness between S26 and S27 could also be hypothesized based on Table III, which shows a correlation of 0.60 between these sessions. However, the classification tree analysis provides us with specific percentages of Entry/NoEntry for S27 based on the values of S26.

The following illustration is for Insight 11, and is given in Table IV. This insight answers the question “*Can the entry of specific entities to an event be predicted?*”. The data mining method used for this purpose is classification analysis. In classification analysis, the dataset is divided into two groups, namely, learning dataset and test dataset. Classification algorithms, also called classifiers (or learners), use the learning dataset to learn from data and predict the class attributes in the test dataset ([59], p17). The prediction success of each learner is measured through classification accuracy (CA) [60], the percentage of correct predictions among all, as well as receiver operating characteristic (ROC) curves [61]. Classifiers which result in higher CA and a greater area under the ROC curve (AUC) correspond to better predictive models. The following classification algorithms are among the best-known classifiers in the machine learning field, and have been used in our analysis: CN2, k-Nearest Neighbor (kNN), Classification Tree, Support Vector Machines (SVM), Naive Bayes, and Neural Networks [59]. Firstly, the entries of entities into event S27 are predicted with a very small learning dataset of 50% (around 130 observations), with 100 experimental repeats (using percentage split of the full dataset into learning and testing datasets). The CA and AUC values are displayed in Table IV, showing that if the behavior of half of the entities for S27 are known, the remaining entry or no entries can be predicted with a very high accuracy, up to 83.15%, with neural network classifier. Besides the black box neural networks technique, which does not tell the reasoning behind classification, CN2 and classification tree might be considered, since they provide the classification rules openly.

D. Behavior Analytics 2: Based on Entity Profile Data

The next illustration is for Insight 12, and is given in Figure 6. This insight answers the question “*Which entities are positioned close to each other?*”. While the question

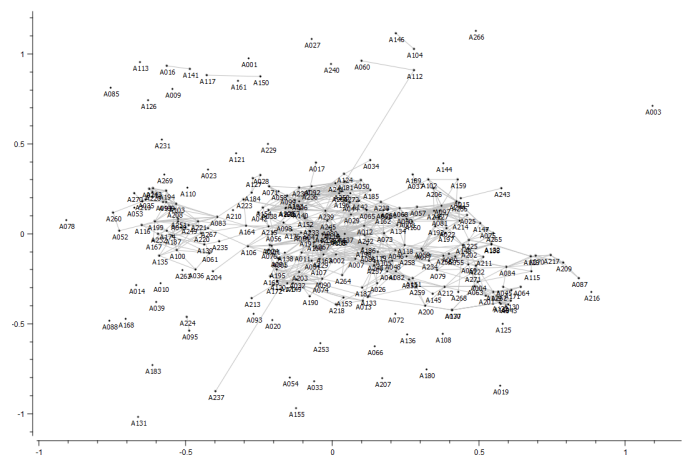


Fig. 6. MDS results for entity analysis based on entity profile data

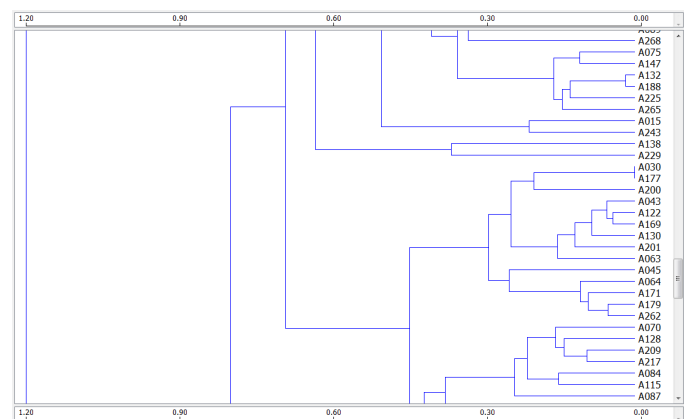


Fig. 7. Hierarchical clustering results for entity analysis based on entity profile data

answered is the same as that of Insight 1, the way it is answered is different. In Insight 1, the answer was computed based on entity-event data, whereas this time it is computed based on entity profile data. The data mining method used for this purpose is again Multi-Dimensional Scaling (MDS). Figure 6 shows the mapping of attendees (entities) on a two-dimensional plane. The most significant associations are shown with lines between the points. Insight 12 is using the database D^6 , which contains only numerical values. Hence, the closeness of the points, as well as the links between them, are based on the Euclidean distance in between. One can observe a highly dense region in Figure 6, to the middle of the figure, as well as a less dense region to the left of the figure, and some sparse links. This means that the entities in the dense cluster are very much close to each other, whereas there are other closely positioned entities among the remaining entities. The results of Insight 12 are different than that of Insight 1, since both the values in the database and the distance measure used are different. This illustrates that one should use the appropriate dataset (and the associated distance measure) that is aligned with the goals of the analysis.

The next illustration is for Insight 13, and is given in Figure 7. This insight answers the question “*Which groups of entities behave most similar?*”. The data mining method

Cluster	Avg_StayDuration	Avg_Earliness	Avg_Lateness	NumberOfEntities
C1	56.31	0.24	-1.45	45
C8	45.97	-15.45	-1.32	23
C4	44.68	-5.65	-5.64	34
C10	43.65	-2.72	28.71	9
C5	41.96	-3.92	-13.05	29
C6	36.76	-12.99	-0.12	12
C3	36.06	-11.68	-15.76	41
C7	32.75	-20.04	-10.70	41
C9	32.18	-6.20	-21.80	17
C2	29.94	-29.59	-5.36	16

Fig. 8. Cluster profiles for entity analysis based on entity profile data

used for this purpose is hierarchical clustering, just as in Insight 6. The data points this time are *entities*, rather than events. So the goal is to see which entities are similar to each other, based on their overall behavior patterns, particularly their entry and exit timings. This analysis uses database \mathcal{D}^6 , and carries hierarchical clustering of entities based on the Euclidean distance between them. The dendrogram in Figure 7 shows that entity groups {A132, A188}, {A030, A177}, {A043, A122, A169}, {A179, A262} are similar to each other, based on the entity profile data.

When partitional clustering is carried out, the entities are partitioned into distinct clusters. One of the analysis to be done given these clusters is to profile the clusters using exploratory data visualization. This cluster profiling constitutes Insight 15, and an illustration of this insight is given in Figure 8. This insight answers the question “*What are the profiles of these entity clusters?*”. Here, the clusters are again based on numerical data coming from the database \mathcal{D}^6 . Figure 8 profiles the clusters based on three attributes, namely average stay duration, average earliness, and average lateness, and also provides the number of entities in each cluster. The 45 entities in the first cluster C1 have the highest average stay duration (56.31 minutes), and have the enter and exit events (sessions in the case study) almost with a perfect timing, neither early nor late. The 23 entities in the next cluster, C8, also stay in the events for a long time (average of 45.97 minutes), and exit with almost no earliness or lateness, but arrive an average of 15.45 minutes late. Each cluster has a profile, that can be similarly read from the figure. For example, at the other extreme, the last cluster, C2, consists of 16 entities who stay the least in the events (average of 29.94 minutes) and enter the events very late (average of 29.59 minutes). A particularly interesting cluster is C10. The 9 entities in cluster C10 stay for a long time in the events, and arrive almost on time, but they stay for a long time (average of 28.71 minutes) after the event is over. In the case study, this can be referring to the after-session discussions participated by these entities.

E. Relationship Network Analysis Based on Entity-Entity Interaction Data

The next illustration is for Insight 16, and is given in Table V. This insight answers the question “*Which entity pairs enter/exit many events together?*”. The data mining method used for this purpose is association mining [62], [63]. The

Attendee1	Attendee2	Support Count
A150	A161	39
A164	A150	31
A009	A150	29
...
A055	A230	6
A249	A230	6
A249	A126	6

TABLE V
THE LIST OF SELECTED ENTITY PAIRS AND THEIR SUPPORT COUNT (ABSOLUTE SUPPORT) VALUES

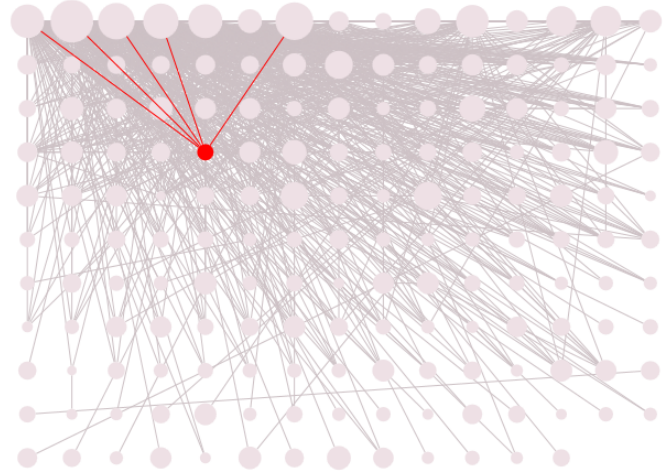


Fig. 9. Results for association mining analysis with grid visualization

database \mathcal{D}^1 is scanned by a pattern matching algorithm (given in Appendix B), and all the entity pairs appearing together are populated into database \mathcal{D}^7 (where an entity pair appears as many times as they are seen together). Then, association mining is carried out to compute the entity pairs that appear together frequently, and this information is populated into database \mathcal{D}^8 . Association mining provides us with the frequent itemsets, namely itemsets that appear together frequently. Only the itemsets that appear at least “minimum support (count) threshold” times are mined and listed. Table V gives a snapshot of \mathcal{D}^8 for our case study, where the minimum support threshold is given in terms of support count (absolute threshold) as 6. So, only the entity pairs that appear together at least 6 times are selected in the association mining analysis and for further analysis. From Table V we can observe that entities {A150, A161} have entered and/or exited together 39 times, which is more than the number of events. This means that they entered and exited together many times during the events, as well, revealing a social connection between these two entities. Other entity pairs with the “strongest” social connection include {A164, A150} and {A009, A150}. It should be noticed here that A150 appears frequently with A161, A164, and A009. So A150 is among the most influential entities. The analysis of “influencer” entities will be extended later in the illustration of Insight 18.

The next illustration is for Insight 17, and is given in Figure 9. This insight answers the question “*How are the entities related to each other?*”. The identified relationship

Node (Attendee)	Degree	Betweenness Centrality
A161	118	2925.52
A150	110	2637.69
A164	99	2470.40
A127	87	1567.21
A009	83	1379.32
A221	71	1012.62
A126	36	445.83
A240	63	359.46
A119	43	350.30
...
A234	1	0.00
A242	1	0.00
A251	1	0.00
A255	1	0.00
A263	1	0.00

TABLE VI
THE LIST OF TOP 10 “INFLUENCERS” AND 5 OF THE “FOLLOWERS”

networks can be used for personalization and generating recommendations for human entities [51]. The data mining methods used for this purpose are network visualization and analysis [64]–[66]. Figure 9 provides a grid visualization of the 163 entities that appear in \mathcal{D}^8 . Each circular node represents an entity; the area of each circle represents the support count (number of times the entity is observed in \mathcal{D}^8); arcs between nodes represent an association between two entities. The visualization is constructed so as to minimize arc crossings. The entities in the upper region of the visualization are those that appear in many interactions. Among those that appear in many interactions, those with smaller area are even more interesting, since we can be inclined in thinking that their interactions were not due to frequent entry-exits, but rather due to interactions with other entities. Furthermore, by visual querying, it is possible to observe how each entity is related to each other entity. In Figure 9, a particular node (entity) is selected and all the associations that it has are highlighted.

Another way of characterizing the nodes (entities) is to compute their graph metrics. One of these metrics is *degree*, which denotes the number of connections for each node. It is an integer value, and it is the summation of in degrees and out degrees of the node. such metric is *betweenness centrality*, which represents the total number of shortest paths for each pair of nodes, if the node is on that path (it can take values between 0 and 1). Detailed information on this and other graph metrics can be found in [65] and [66].

The next illustration is for Insight 18, and is given in Table VI. This insight answers the question “Which entities are influencers and which are followers?”. The data mining method used for this purpose is network analysis, and specifically network characterization through computing node metrics [65], [66]. Table VI lists the most active and influential entities (A161 through A119), as well as some of the least associated ones (which appear only 6 times with other entities). The most influential entities have high value for betweenness centrality, indicating that they are at the “crossroads” of social networks.

The next illustration is for Insight 19, and is given in Figure 10. This insight answers the question “How can the behavioral

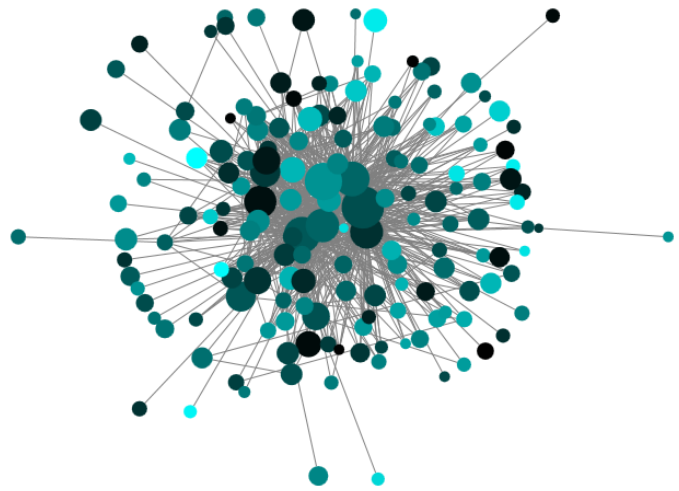


Fig. 10. Results for association mining analysis with Harel-Koren layout algorithm

attributes be analyzed together with the social network?”. This analysis is specifically aimed at the scenarios where the entities are humans. The data mining method used for this purpose is network visualization, where some behavioral attributes are mapped onto the nodes. The network in Figure 10 is exactly the same as that in Figure 9 with respect to node-link structure. However, the selected layout algorithm is different (Harel-Koren algorithm; [67]), and nodes are colored according to average stay duration (a behavioral attribute). Lighter colors denote longer average stay durations. Size again denotes the support count of the node. The visualization in Figure 10 is constructed using \mathcal{D}^9 , which is obtained by joining the two different databases of \mathcal{D}^6 (entity profile data) and \mathcal{D}^8 (association graph). The nodes in the center are influencers and the ones on the outside are followers. However, we are now also able to see which influencers stay in the events for long, and which stay less. Thus, we are able to see the influencers that enhance our desired goals (longer stay durations in sessions with less frequent entry-exits, in our case study) versus the influencers that disrupt the system (by staying very short in sessions and entering and exiting many times). So we are not only able to identify the social network and the influence the entities have on the network, but also the direction of the effects (positive or negative) that they have.

V. PRACTICAL IMPLICATIONS

In this section, we discuss how the insights and information obtained through our analysis can be used in a multitude of ways, for improving the system and achieving various goals.

First, information regarding similar-behaving entities in a schedule-based system can be used in several ways:

- In the context of social event management, ubiquitous information systems can use this information to suggest new people for professional social networks. For example, in a conference, when two attendees are identified as entering and exiting similar events, the conference mobile application can recommend them each other to add into LinkedIn and other professional social networks. Another

1 use of the information is the suggestion of events to social
 2 event attendees in which similar-behaving attendees have
 3 already entered.

- 4 • In the context of manufacturing, if one of the entities
 5 has entered the production system, a-priori planning can
 6 be done for similar-behaving entities. Furthermore, the
 7 production system can be set up to accommodate not
 8 only the already entered entity, but also those that may
 9 potentially enter, in a way to reduce total setup time.
- 10 • In the context of warehousing, an example scenario
 11 where the information on similar-behaving entities can
 12 be used is the following: Consider pallets of similar-
 13 behaving products entering the warehouse. There is a high
 14 chance that they will also exit together. Therefore, the
 15 warehouse management system (WMS) software can be
 16 programmed so as to allocate neighboring locations for
 17 these two pallets, so that they can be put away and picked
 18 on the same route, saving time and cost.
- 19 • In the context of healthcare, similar-behaving entities can
 20 be equipment used for surgical operations. In this sce-
 21 nario, these equipment can be stored in the same storage
 22 room when not in use. This way, they can be accessed in
 23 the least possible time when an urgent surgical operation
 24 is to be conducted.
- 25 • In the context of education, similar-behaving students can
 26 be identified automatically based on their entries and exits
 27 to classrooms. Then this information is populated into the
 28 school's information system databases. In case a student
 29 can't be reached by phone, the school management or
 30 the instructors can try to reach him through contacting
 31 his friend.
- 32 • Finally, in the context of tourism, visitors in a museum
 33 can be offered special places of interest in the museum
 34 (visited by similar-behaving visitors) through the smart
 35 mobile devices that are guiding them.

36 Information regarding groups of similar events in a
 37 schedule-based system can also be used in a multitude of
 38 ways for improving the system and achieving various goals.
 39 One example use case from social event management is the
 40 joint design and improvement of similar sessions in the social
 41 event. The session managers can come together and discuss
 42 possible opportunities of improving the similar sessions in
 43 future conferences.

44 Information regarding the correlation between events can
 45 be used for improving schedules. For example, in the context
 46 of manufacturing, successive production periods with negative
 47 correlation may experience great changes in the product mix
 48 entering the production. These can also be sources of long
 49 setup times and costs. Therefore, schedule can be adjusted
 50 based on the results of data mining.

51 Information regarding earlier events affecting later events,
 52 as well the predictability of entry of specific entities to an
 53 event, can also be used for benefiting the system. Consider a
 54 warehousing scenario where the entities are the various types
 55 of products loaded on pallets. Based on the past behaviors
 56 of these products, one can estimate for each product the
 57 probability of entering a particular warehouse zone during a
 58 particular time interval (event). This can be used to predict
 59
 60

whether capacity will be exceeded in that zone of the ware-
 house during that time interval. Then, if necessary, additional
 capacity can be created, for example, through establishing
 temporary additions to that zone using pallet stacking frames.

Finally, the insights regarding influencers in the system
 can be utilized in many ways. For example, in social event
 management, once these influencers are determined, they can
 be consulted for help in promoting newly established sessions
 or for increasing membership to the organizing society.

VI. CONCLUSIONS AND FUTURE RESEARCH

The importance of RFID systems for data collection and
 processing is ever increasing. RFID systems find applications
 in a very wide range of domains, including in schedule-
 based systems, which operate based on (or contain within)
 a schedule of events. In this paper, we have presented a
 comprehensive framework for mining of RFID data coming
 from schedule-based systems, for the first time in the literature.
 Our framework is generic, and can be applied to any schedule-
 based system that operates as described.

There exists two very fundamental future research avenues
 for extending the current work:

- RFID tags can collect and/or carry not only location
 and time information, but other information, as well.
 Such information typically includes entity type, entity
 affiliation, physical attributes, and assigned attributes. In a
 logistics context, examples of these attributes are product
 type, manufacturer, weight, and price [52]. The additional
 information may also be collected through various sen-
 sors (e.g. temperature, GPS) integrated within or mounted
 on the tags. While the framework that we have presented
 here considers only time data in relation to schedule data,
 it can be highly enriched with the incorporation of analy-
 sis of these additional attributes. For example, scheduling
 and plans are important in manufacturing context, and
 are very much dependent on quality level achieved in
 the production process. Quality-related attributes can be
 analyzed together with product attributes and schedule
 data to improve the production process in the dimensions
 of quality, time, and cost. To this end, the re-mining
 framework of [68] can be integrated with the framework
 here to augment the data and to discover further insights.
- The other fundamental research avenue is extending the
 framework from the temporal domain to the spatio-
 temporal domain, by extending it to handle multiple
 locations.
- While the analysis of serial events is fundamental, con-
 sideration can be made in future research for concu-
 rrent events (events that can independently take place
 at the same time) in the system. The consideration for
 concurrent events would require significant changes in
 the augmentation algorithm, as it would require complex
 event processing [17]. However, the applicability of the
 current framework, as well as the types of analysis and
 the insights obtained, would still be relevant and useful.
- One of the important challenges in industrial applications
 is the challenge of big data [69]. A possible future

research can involve the development of the framework to accommodate for big data applications. To support the large volumes of input data, when the proposed framework is implemented, the data processing of this framework should be split into independent tasks to support parallel processing systems such as MapReduce. As indicated by Figure 2 in the manuscript, our framework has very few interactions between different branches of data flows and thus splitting the overall data processing into multiple tasks is possible and can be highly feasible. The methods for MapReduce implementation of the individual data mining and data visualization algorithms used in this manuscript, such as hierarchical clustering, can be found in the literature [70]. Hence, the proposed framework can support the big data environment if its implementation is properly designed.

Other possible future research avenues include the following:

- The concepts and methods used for the analysis of behavior in electronic games and virtual worlds [49] [71] [72] can be used in the analysis of RFID data, and vice versa.
- The methods used for analyzing animal societies based on RFID data [73] can be adopted to analyzing the movement of entities in schedule-based systems in general.
- Data from RFID (and other types of sensors) have been used in some literature [4], [74], [75] to (optimally) allocate the RFID readers. Data mining frameworks can be integrated with such methods to come up with better allocation of reader within an environment.
- The study can be extended such that it encompasses more of the available data mining algorithms and techniques. For example, besides using k-Means Clustering in the unsupervised learning process, one can use k-Means++ [76], to reduce both clustering errors and running times.
- Last but not least, mining of RFID data can be used in the general context of ambient intelligence applications, which are surveyed and discussed in [77]–[80].

APPENDIX A. AUGMENTATION ALGORITHMS

The first augmentation algorithm is the following:

Input: $\mathcal{D}^0, \mathcal{D}'$

Output: \mathcal{D}^1

```

foreach  $d^0 \in \mathcal{D}^0 = \langle r, i, u_{ir}, U_{ir} \rangle$  do
   $d^1 = \langle r, i, u_{ir}, U_{ir}, j_1 = intervalOf(u_{ir}),$ 
     $j_2 = intervalOf(U_{ir}), \varepsilon_1 = intervalType(j_1),$ 
     $\varepsilon_2 = intervalType(j_2), s_{j_1}, s_{j_2}, f_{j_1}, f_{j_2} \rangle;$ 
   $\mathcal{D}^1 \sqcup d^1;$ 
end

```

Algorithm 1: Generate \mathcal{D}^1

By using the schedule data \mathcal{D}' , this algorithm augments each record of the RFID database \mathcal{D}^0 with the information of the intervals covering the entry time and exit time. This information includes the type of interval, start time, and finish time. The augmented records forms a new database \mathcal{D}^1 for further analysis. Algorithm 1 includes a single loop that requires the initial construction of the lookup tables for the lookup

functions. Each lookup has to scan through the J intervals for each of the R records. Running time of this initialization stage is $O(RJ)$. After this, each record is augmented, taking $O(R)$ time. So the running time of Algorithm 1 is $O(RJ)$.

The second augmentation algorithm is given below.

For each record of database \mathcal{D}^1 , this algorithm first identifies the types of the sequence of intervals that partially or completely falls between the entry time u and exit time U of that particular record. If an interval j is an event, its entry time scenario is analyzed to derive e and p . This is followed by the analysis of exit time scenario to determine x and q . Finally, the time T spent on that event, the *earliness'* and the *lateness'* are computed for that event based on $e, p, x,$ and q . The first four fields in this record are then augmented with intermediary data e, p, x, q as well as the computed $t, \text{earliness}'$ and $\text{lateness}'$. The augmented new record is added to a new database \mathcal{D}^2 . This procedure is repeated for all records in database \mathcal{D}^1 . Algorithm 2 has two interleaved loops, and runs for each record and for each interval. So the running time of Algorithm 2 is $O(RJ)$.

In the below algorithm, by noting that two successive break intervals are impossible and at least part of the event falls within the time span bounded by u and U , the following possible entry time scenarios for the event are considered:

- 1) Entry time u falls within the event
- 2) Entry time u is before the start time of the event
 - a) The interval $j - 1$ preceding the event and in the sequence is a break
 - i) The second interval $j - 2$ preceding the event and in the sequence is an event
 - ii) The second interval $j - 2$ preceding the event and in the sequence does not exist (Type of interval is *null*)
 - b) The interval $j - 1$ preceding the event and in the sequence is an event

```

1  Input:  $\mathcal{D}^1$ 
2  Output:  $\mathcal{D}^2$ 
3  foreach  $d^1 \in \mathcal{D}^1$  do
4       $u = d^1 \cdot u_{ir}$ ;
5       $U = d^1 \cdot U_{ir}$ ;
6      foreach  $j \in \{j : j_1 \leq j \leq j_2\}$  do
7          if  $intervalType(j) = Event$  then
8              /* Analyze different entry time scenarios */
9              if  $u \geq s_j$  then
10                  $e = u$ ;
11                  $p = u$ ;
12                  $entryStatus = LateEntry$ ;
13             else
14                  $e = s_j$ ;
15                 if  $intervalType(j-1) = Break$  then
16                      $entryStatus = EarlyEntry$ ;
17                     if  $intervalType(j-2) = Event$  then
18                          $p = s_j - (f_{j-1} - s_{j-1}) / 2$ ;
19                     end
20                     if  $intervalType(j-2) = null$  then
21                          $p = u$ ;
22                     end
23                 end
24                 if  $intervalType(j-1) = Event$  then
25                      $entryStatus = EntryFromPreEvent$ ;
26                      $p = s_j$ ;
27                 end
28             end
29             /* Analyze different exit time scenarios */
30             if  $U \leq f_j$  then
31                  $x = U$ ;
32                  $q = U$ ;
33                  $exitStatus = EarlyExit$ ;
34             else
35                  $x = f_j$ ;
36                 if  $intervalType(j+1) = Break$  then
37                      $exitStatus = LateStatus$ ;
38                     if  $intervalType(j+2) = Event$  then
39                          $q = f_j + (f_{j+1} - s_{j+1}) / 2$ ;
40                     end
41                     if  $intervalType(j+2) = null$  then
42                          $q = U$ ;
43                     end
44                 end
45                 if  $intervalType(j+1) = Event$  then
46                      $exitStatus = ExitIntoNextEvent$ ;
47                      $q = f_j$ ;
48                 end
49             end
50             /* Compute the metrics */
51              $T = x - e$ ;
52              $earliness' = s_j - p$ ;
53              $lateness' = q - f_j$ ;
54         end
55          $d \leftarrow \langle i, j, r, u, U, e, x, T, p, q, earliness', lateness', entryStatus, exitStatus \rangle$ ;
56          $\mathcal{D}^2 \sqcup d$ ;
57     end
58 end

```

Algorithm 2: Generate \mathcal{D}^2

Similarly, the possible exit time scenarios under consideration are

- 1) Exit time U falls within the event
- 2) Exit time U is after the finish time of the event
 - a) The interval $j + 1$ following the event and in the sequence is a break
 - i) The second interval $j + 2$ following the event and in the sequence is an event
 - ii) The second interval $j + 2$ following the event and in the sequence does not exist (Type of interval is *null*)
 - b) The interval $j + 1$ following the event and in the sequence is an event

Different scenarios of entry time and exit time have different expressions for e , x , p , q respectively as presented in the algorithm. The computation for the metrics, however, is unified in the sense that it is independent of the types of scenarios.

APPENDIX B. PATTERN MATCHING ALGORITHM

The pattern matching algorithm is given below:

```

Input:  $\mathcal{D}^1$ 
Output:  $\mathcal{D}^{7a}$  and  $\mathcal{D}^{7b}$ 
 $\mathcal{N}^1$ : Number of records in  $\mathcal{D}^1$ 
 $\mathcal{T}_{RD}$ : Time condition for relationship detection
for  $m \leftarrow 1$  to  $\mathcal{N}^1 - 1$  do
     $d_m^1 = \text{the } m^{\text{th}} \text{ record in } \mathcal{D}^1$ ;
    for  $k \leftarrow m + 1$  to  $\mathcal{N}^1$  do
         $d_k^1 = \text{the } k^{\text{th}} \text{ record in } \mathcal{D}^1$ ;
        if  $d_m^1 \cdot i \neq d_k^1 \cdot i$  then
            if  $|d_m^1 \cdot u_{ir} - d_k^1 \cdot u_{ir}| \leq \mathcal{T}_{RD}$  then
                 $\mathcal{D}^{7a} \sqcup \langle d_m^1 \cdot i, d_k^1 \cdot i \rangle$ ;
            end
            if  $|d_m^1 \cdot U_{ir} - d_k^1 \cdot U_{ir}| \leq \mathcal{T}_{RD}$  then
                 $\mathcal{D}^{7b} \sqcup \langle d_m^1 \cdot i, d_k^1 \cdot i \rangle$ ;
            end
        end
    end
end

```

Algorithm 3: Generate \mathcal{D}^{7a} and \mathcal{D}^{7b} whose union forms \mathcal{D}^7

This algorithm first generates all possible combinations of two arbitrary records from \mathcal{D}^1 that have different entities. For each combination, if the entry time difference between the two records is less than or equal to a predefined length \mathcal{T}_{RD} , the transaction ID for entry time is first updated. Then two new records with updated transaction ID, entity, entry time and record ID are generated and added to a database \mathcal{D}^{7a} . Similarly, if the exit time difference is less than or equal to \mathcal{T}_{RD} , the transaction ID for exit time is updated and two new records are inserted into a database \mathcal{D}^{7b} . The union of \mathcal{D}^{7a} and \mathcal{D}^{7b} form a new database \mathcal{D}^7 containing all detected pairs showing the close relationship between entry or exit time of two entities for a particular event. Algorithm 3 has two interleaved loops, each executed for up to R records. So the running time of Algorithm 3 is $O(R^2)$.

APPENDIX C. DATA MINING PROCESSES

Figures 11 and 12 display the data mining processes carried out. The first process in Figure 11 shows an unsupervised machine learning model, whereas the second process shown in Figure 12 shows a supervised machine learning model.

The unsupervised data mining process (Figure 11) starts with reading data from file (File block), verifying that the data is read correctly (Data Table 1 block), and handling any missing values (Impute block). Next, data is again verified, this time visually, using a scatter plot (Scatterplot 1 block). The attributes are selected and specified (Select Attributes) and the unsupervised learning is initiated. The first type of analysis uses entity-entity distances (Example Distance) and conducts MDS (MDS), as well as Hierarchical Clustering, and detects any Outliers. The next analysis is k-Means Clustering, whose results are visually inspected (Scatterplot 2) and exported into a data table (Data Table 2). The final analysis is the computation of distances between events (Attribute Distance) and the conduct of hierarchical clustering (Hierarchical Clustering 2).

The supervised data mining process (Figure 12) also starts with the same steps. However, the attribute selection is different, because -unlike the previous process- one categorical attribute (S27, in our case study) has to be selected as the class attribute to be predicted. Next, multiple classifiers are

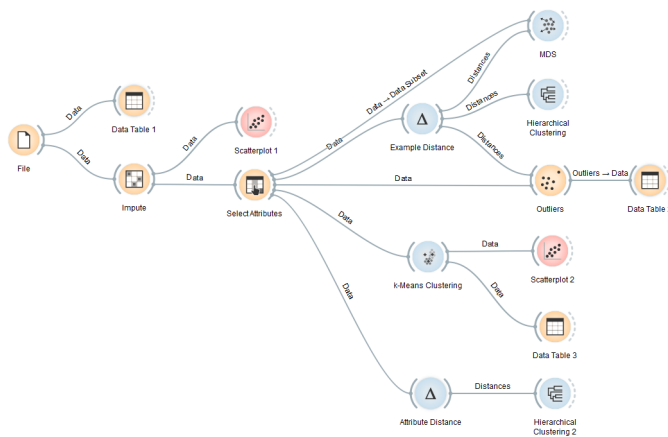


Fig. 11. The unsupervised data mining process conducted in the study

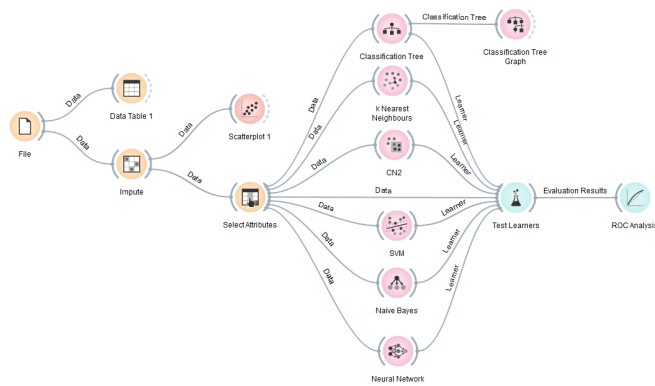


Fig. 12. The supervised data mining process conducted in the study

applied, and their performances are tested and compared (Test Learners and ROC Analysis blocks). One of the classifiers has an additional benefit. Besides being used in classification analysis, the Classification Tree classifier is used in constructing the Classification Tree Graph.

APPENDIX D. SOFTWARE TOOLS USED

There exist a multitude of data analysis and data mining software tools, and we have used different tools for different purposes. Matlab¹ was used for coding the developed and presented algorithms. Orange² data mining software [81] was used for clustering, classification, and classification tree analysis. RapidMiner³ [82] was used to compute the correlation matrix for the sessions. Borgelt's implementation of the apriori algorithm⁴ [83]–[85] was used to compute frequent itemsets (attendees frequently appearing together). Finally, NodeXL⁵ [86] was used to visualize association mining results and to compute graph metrics, enabling association-based social network analysis.

¹<http://www.mathworks.com>

²<http://orange.biolab.si/>

³<http://rapidminer.com/>

⁴<http://www.borgelt.net/apriori.html>

⁵<http://nodexl.codeplex.com/>

ACKNOWLEDGEMENTS

The authors thank Akın Altunbaş and Ahmet Erdem Altunbaş from Borda Technology and Enes Eryarsoy from İstanbul Şehir University for introducing the problem to the research group and providing the data for the case study. The authors also thank Utku Kaymaz, Berkay Dönmez, and Çağrı Başel from Sabancı University for their assistance in the editing of the paper.

REFERENCES

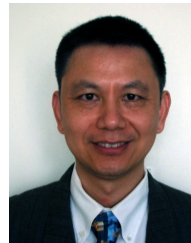
- [1] R.J. Schonberger. Applications of single-card and dual-card kanban. *Interfaces*, 13(4), 56-67, 1983.
- [2] W.C. Benton. Push and Pull Production Systems. *Wiley Encyclopedia of Operations Research and Management Science*. Wiley, 2011.
- [3] X. Zhu, S.K. Mukhopadhyay, H. Kurata. A review of RFID technology and its managerial applications in different industries. *Journal of Engineering and Technology Management*, 29(1), 152-167, 2012.
- [4] A. Oztekin, F.M. Pajouh, D. Delen, L.K. Swim. An RFID network design methodology for asset tracking in healthcare. *Decision Support Systems*, 49(1), 100-109, 2010.
- [5] W.-P. Liao, T.M.Y. Lin, S.-H. Liao. Contributions to Radio Frequency Identification (RFID) research: An assessment of SCI-, SSCI-indexed papers from 2004 to 2008. *Decision Support Systems*, 50, 548-556, 2011.
- [6] J. Han, H. Gonzalez, X. Li, D. Klabjan. Warehousing and mining massive RFID data sets. In *Advanced Data Mining and Applications* (pp. 1-18). Springer Berlin Heidelberg, 2006.
- [7] E.N. Cincioğlu, P.P. Shenoy, C. Kocabasoglu. Use of radio frequency identification for targeted advertising: a collaborative filtering approach using Bayesian networks. In *Symbolic and Quantitative Approaches to Reasoning with Uncertainty* (pp. 889-900). Springer Berlin Heidelberg, 2007.
- [8] J. Han, M. Kamber, J. Pei. *Data Mining: Concepts and Techniques*. Morgan Kaufmann. 3 edition. Burlington, MA, USA, 2011.
- [9] M.L. Pinedo. *Scheduling: Theory, Algorithms, and Systems*. Prentice Hall College Div, 1st ed., 1994.
- [10] M.L. Pinedo. *Planning and Scheduling in Manufacturing and Services*. Springer, 2nd ed., 2009.
- [11] Y. Yin, M. Liu, J. Hao, M. Zhou. Single-Machine Scheduling With Job-Position-Dependent Learning and Time-Dependent Deterioration. *IEEE Transactions on Systems, Man, and Cybernetics - Part A: Systems and Humans*, 42(1), 192-200, 2012.
- [12] J.S.K. Lau, G.Q. Huang, K.L. Mak, L. Liang. Agent-Based Modeling of Supply Chains for Distributed Scheduling. *IEEE Transactions on Systems, Man, and Cybernetics - Part A: Systems and Humans*, 36(5), 847-861, 2006.
- [13] X. Qiu, H.Y.K. Lau. An AIS-based hybrid algorithm for static job shop scheduling problem. *Journal of Intelligent Manufacturing*, 25(3), 489-503, 2014.
- [14] R. Balasundaram, N. Baskar, R. Siva Sankar. Discovering dispatching rules for job shop scheduling using data mining. *Advances in Intelligent Systems and Computing*, 178, 63-72, Springer Berlin Heidelberg, 2013.
- [15] C. Rainer. Data Mining as Technique to Generate Planning Rules for Manufacturing Control in a Complex Production System A Case Study from a Manufacturer of Aluminum Products. In: K. Windt (ed.), *Robust Manufacturing Control, Lecture Notes in Production Engineering*, 2013.
- [16] C.L. Wang, G. Rong, W. Weng, Y.P. Feng. Mining scheduling knowledge for job shop scheduling problem. *15th IFAC Symposium on Information Control Problems in Manufacturing ? INCOM 2015* 48(3), 800-805, 2015.
- [17] R. Helouai, M. Niepert, H. Stuckenschmidt. Recognizing interleaved and concurrent activities using qualitative and quantitative temporal relationships. *Pervasive and Mobile Computing*, 7, 660-670, 2011.
- [18] G. Mariscal, O. Marban, C. Fernandez. A survey of data mining and knowledge discovery process models and methodologies. *The Knowledge Engineering Review*, 25(02), 137-166, 2010.
- [19] S. Sharma, K.M. Osei-Bryson, G.M. Kasper. Evaluation of an integrated Knowledge Discovery and Data Mining process model. *Expert Systems with Applications*, 39(13), 11335-11348, 2012.
- [20] W.S. Ku, H. Chen, H. Wang, M.T. Sun. A Bayesian Inference-Based Framework for RFID Data Cleansing. *IEEE Transactions on Knowledge and Data Engineering*. doi:10.1109/TKDE.2012.116, 2012.

- [21] A.I. Baba, H. Lu, X. Xie, T.B. Pedersen. Spatiotemporal data cleansing for indoor RFID tracking data. In *Mobile Data Management (MDM)*, 2013 IEEE 14th International Conference on (Vol. 1, pp. 187-196). IEEE, 2013.
- [22] T.C. Poon, K.L. Choy, H.K. Chow, H.C. Lau, F.T. Chan, K.C. Ho. A RFID case-based logistics resource management system for managing order-picking operations in warehouses. *Expert Systems with Applications*, 36(4), 8277-8301, 2009.
- [23] A. Ilic, T. Andersen, F. Michahelles. Increasing supply-chain visibility with rule-based RFID data analysis. *Internet Computing, IEEE*, 13(1), 31-38, 2009.
- [24] D. Shuping, W. Wright. Geotime Visualization of RFID. *RFID J.*, Mar./Apr. 2005, 1-6, 2005.
- [25] S. Miyazaki, T. Washio, K. Yada. Analysis of Residence Time in Shopping Using RFID Data - An Application of the Kernel Density Estimation to RFID. In *Data Mining Workshops (ICDMW)*, 2011 IEEE 11th International Conference on (pp. 1170-1176). IEEE, 2011.
- [26] B. Fang, S. Liao, K. Xu, H. Cheng, C. Zhu, H. Chen. A novel mobile recommender system for indoor shopping. *Expert Systems with Applications*, 39(15), 11992-12000, 2012.
- [27] S. Sakurai, M. Sanbe, K. Watanabe. Application of the RFID Data Mining to an Apparel Field. In *Network-Based Information Systems (NBIS)*, 2010 13th International Conference on (pp. 28-35). IEEE, 2010.
- [28] J. Lyu Jr, S.Y. Chang, T.L. Chen. Integrating RFID with quality assurance system - Framework and applications. *Expert Systems with Applications*, 36(8), 10877-10882, 2009.
- [29] C.K.H. Lee, K.L. Choy, G.T. Ho, K.M.Y. Law. A RFID-based Resource Allocation System for garment manufacturing. *Expert Systems with Applications*, 40, 784-799, 2012.
- [30] W. Wen. An intelligent traffic management expert system with RFID technology. *Expert Systems with Applications*, 37(4), 3024-3035, 2010.
- [31] C.Y. Tsai, J.J. Liou, C.J. Chen, C.C. Hsiao. Generating touring path suggestions using time-interval sequential pattern mining. *Expert Systems with Applications*, 39(3), 3593-3602, 2012.
- [32] Y. Meiller, S. Bureau, W. Zhou, S. Piramuthu. Adaptive knowledge-based system for health care applications with RFID-generated information. *Decision Support Systems*, 51(1), 198-207, 2011.
- [33] J. Lapalu, K. Bouchard, A. Bouzouane, B. Bouchard, S. Giroux. Unsupervised Mining of Activities for Smart Home Prediction. *Procedia Computer Science*, 19, 503-510, 2013.
- [34] H.H. Hsu, Z. Cheng, T.K. Shih, C.C. Chen. RFID-Based Personalized Behavior Modeling. In *Ubiquitous, Autonomic and Trusted Computing*, 2009. UIC-ATC'09. Symposia and Workshops on (pp. 350-355). IEEE, 2009.
- [35] H.H. Hsu, C.C. Chen. RFID-based human behavior modeling and anomaly detection for elderly care. *Mobile Information Systems*, 6(4), 341-354, 2010.
- [36] M. Delgado, M. Ros, M. Amparo Vila. Correct behavior identification system in a Tagged World. *Expert Systems with Applications*, 36(6), 9899-9906, 2009.
- [37] Y. Liu, L. Chen, J. Pei, Q. Chen, Y. Zhao. Mining frequent trajectory patterns for activity monitoring using radio frequency tag arrays. In *Pervasive Computing and Communications*, 2007. PerCom'07. Fifth Annual IEEE International Conference on (pp. 37-46). IEEE, 2007.
- [38] E. Masciari. A Framework for Outlier Mining in RFID data. In *Database Engineering and Applications Symposium*, 2007. IDEAS 2007. 11th International (pp. 263-267). IEEE, 2007.
- [39] C. Cattuto, W. Van den Broeck, A. Barrat, V. Colizza, J.F. Pinton, A. Vespignani. Dynamics of person-to-person interactions from distributed RFID sensor networks. *PLoS one*, 5(7), e11596, 2010.
- [40] H. Gao, H. Liu. Data Analysis on Location-Based Social Networks. In: Chin, A., Zhang, D., *Mobile Social Networking: An Innovative Approach*. Springer, 2014.
- [41] T.-S. Chen, Y.-S. Chou, T.-C. Chen. Mining User Movement Behavior Patterns in a Mobile Service Environment. *IEEE Transactions on Systems, Man, and Cybernetics - Part A: Systems and Humans*, 42(1), 87-101, 2012.
- [42] P. Kazienko, K. Musial, T. Kajdanowicz. Multidimensional Social Network in the Social Recommender System. *IEEE Transactions on Systems, Man, and Cybernetics - Part A: Systems and Humans*, 41(4), 746-759, 2011.
- [43] M. Szomszor, P. Kostkova, C. Cattuto, W. Van den Broeck, A. Barrat, H. Alani. Providing Enhanced Social Interaction Services for Industry Exhibitors at Large Medical Conferences. In *Developments in E-systems Engineering (DeSE)*, 2011 (pp. 42-45). IEEE, 2011.
- [44] A. Chin, B. Xu, F. Yin, X. Wang, W. Wang, et al. Using proximity and homophily to connect conference attendees in a mobile social network. In *Distributed Computing Systems Workshops (ICDCSW)*, 2012 32nd International Conference on (pp. 79-87). IEEE, 2012.
- [45] W. Reinhardt, T. Messerschmidt, T. Nelkner. Awareness-support in scientific events with SETapp. In *Proceedings of the 1st European Workshop on Awareness and Reflection in Learning Networks*, 2011.
- [46] J. Bravo, R. Hervás, I. Sánchez, G. Chavira, S.W. Nava. Visualization Services in a Conference Context: An Approach by RFID Technology. *J. UCS*, 12(3), 270-283, 2006.
- [47] I. Hsu. Extending UML to model Web 2.0-based context-aware applications. *Software: Practice and Experience*, 42(10), 1211-1227, 2012.
- [48] M. Atzmueller. Mining social media: key players, sentiments, and communities. *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery*, 2, (5), 411-419, 2012.
- [49] D.A. Huffaker, C. Teng, M.P. Simmons, L. Gong, L.A. Adamic. Group Membership and Diffusion in Virtual Worlds. In *Privacy, security, risk and trust*, 2011 IEEE Third International Conference on Social Computing (SOCIALCOM) (pp. 331-338). IEEE, 2011.
- [50] J.J. Jung. Ubiquitous conference management system for mobile recommendation services based on mobilizing social networks: A case study of u-conference. *Expert Systems with Applications*, 38(10), 12786-12790, 2011.
- [51] M. Atzmueller, S. Doerfel, A. Hotho, F. Mitzlaff, G. Stumme. Face-to-face contacts during a conference: Communities, roles, and key players. In *The Second International Workshop on Mining Ubiquitous and Social Environments* (pp. 25), 2011.
- [52] H. Gonzalez, J. Han, H. Cheng, X. Li, D. Klabjan, T. Wu. Modeling massive RFID data sets: a gateway-based movement graph approach. *IEEE Transactions on Knowledge and Data Engineering*, 22(1), 90-104, 2010.
- [53] Y. Wang, E.P. Lim, S.Y. Hwang. Efficient mining of group patterns from user movement data. *Data & Knowledge Engineering*, 57(3), 240-282, 2006.
- [54] I. Borg, P.J.F. Groenen, P. Mair. *Applied Multidimensional Scaling*. Springer, 2012.
- [55] R. Bose. *Information Theory, Coding and Cryptography*. page 102. 2008.
- [56] L. Kaufman, P.J. Rousseeuw. *Finding groups in data: an introduction to cluster analysis*. Wiley Series in Probability and Statistics (Book 603). Wiley-Interscience; 1st edition, 2005.
- [57] J.L. Rodgers, W.A. Nicewander. Thirteen ways to look at the correlation coefficient. *The American Statistician*, 42(1):59-66, February 1988, 2009.
- [58] L. Rokach, O. Maimon. *Data mining with decision trees: theory and applications*. World Scientific Pub Co Inc, 2008.
- [59] E. Alpaydin. *Introduction to Machine Learning*, The MIT Press, Cambridge, MA, 2010.
- [60] C.E. Brodley, M.A. Friedl. Identifying Mislabeled Training Data. *Journal of Artificial Intelligence Research*, 11, 131-167, 1999.
- [61] T. Fawcett. ROC Graphs: Notes and Practical Considerations for Researchers, *Pattern Recognition Letters*, 27(8), 882-891, 2004.
- [62] R. Agrawal, T. Imielinski, A. Swami. Mining association rules between sets of items in large databases. In *ACM SIGMOD Record* (Vol. 22, No. 2, pp. 207-216). ACM, 1993.
- [63] Agrawal, R., & R. Srikant. Fast algorithms for mining association rules. In *Proc. 20th Int. Conf. Very Large Data Bases, VLDB* (Vol. 1215, pp. 487-499), 1994.
- [64] I. Herman, G. Melancon, M.S. Marshall. Graph visualization and navigation in information visualization: A survey. *IEEE Transactions on Visualization and Computer Graphics*, 6(1), 24-43, 2000.
- [65] C. Christensen, A. Réka. Using graph concepts to understand the organization of complex systems. *International Journal of Bifurcation and Chaos*, 17, 2201-2214, 2007.
- [66] T. Opsahl, F. Agneessens, J. Skvoretz. Node centrality in weighted networks: Generalizing degree and shortest paths. *Social Networks*, 32, 245, 2010.
- [67] D. Harel, Y. Koren. Graph drawing by high-dimensional embedding. In *Graph Drawing* (pp. 207-219). Springer Berlin Heidelberg, 2002.
- [68] A. Demiriz, G. Ertek, T. Atan, U. Kula. Re-mining item associations: Methodology and a case study in apparel retailing. *Decision Support Systems*, 52(1), 284-293, 2011.
- [69] V. Mayer-Schönberger, K. Cukier. *Big Data: A Revolution That Will Transform How We Live, Work, and Think*. Eamon Dolan/Mariner Books, 2014.
- [70] T. Sun, C. Shuy, F. Liy, H. Yuy, L. Ma, Y. Fang. An efficient hierarchical clustering method for large datasets with Map-Reduce. In *PDCAT*, 2009.
- [71] K. Borner, S. Penumathy. Social diffusion patterns in three-dimensional virtual worlds. *Information Visualization*, 2(3), 182-198, 2003.

- 1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
- [72] N. Hoobler, G. Humphreys, M. Agrawala. Visualizing competitive behaviors in multi-user virtual environments. In Proceedings of the conference on Visualization'04 (pp. 163-170). IEEE Computer Society, 2004.
- [73] G. Cabanes, Y. Bennani, D. Fresneau. Mining RFID behavior data using unsupervised learning. *International Journal of Applied Logistics*, 1(1), 28-47, 2010.
- [74] W. Chang, D. Zeng, H. Chen. A stack-based prospective spatio-temporal data analysis approach. *Decision Support Systems*, 45(4), 697-713, 2008.
- [75] B.C. Cheung, S.L. Ting, A.H. Tsang, W.B. Lee. A methodological approach to optimizing RFID deployment. *Information Systems Frontiers*, 1-15, 2012.
- [76] D. Arthur, S. Vassilvitskii. k-means++: The advantages of careful seeding. *Proceedings of the Eighteenth Annual ACM-SIAM Symposium on Discrete Algorithms*, 1027-1035, Society for Industrial and Applied Mathematics, 2007.
- [77] C. Ramos, J.C. Augusto, D. Shapiro. Ambient intelligence-The next step for artificial intelligence. *Intelligent Systems, IEEE*, 23(2), 15-18, 2008.
- [78] D.J. Cook, J.C. Augusto, V.R. Jakkula. Ambient intelligence: Technologies, applications, and opportunities. *Pervasive and Mobile Computing*, 5(4), 277-298, 2009.
- [79] F. Sadri. Ambient intelligence: A survey. *ACM Computing Surveys (CSUR)*, 43(4), 36, 2011.
- [80] J.C. Augusto, H. Nakashima, H. Aghajan. Ambient intelligence and smart environments: A state of the art. In *Handbook of Ambient Intelligence and Smart Environments* (pp. 3-31). Springer US, 2010.
- [81] J. Demšar, T. Curk, A. Erjavec, Č Gorup, T. Hočevcar, M. Milutinovič, M. Možina, M. Polajnar, M. Toplak, A. Starič, M. Štajdohar, L. Umek, L. Žagar, J. Žbontar, M Žitnik, B. Zupan. Orange: data mining toolbox in Python. *Journal of Machine Learning Research*, 14, 2349-2353, 2013.
- [82] M. Hofmann, R. Klinkenberg. *RapidMiner: Data Mining Use Cases and Business Analytics Applications*. CRC Press, 2013.
- [83] C. Borgelt, R. Kruse. Induction of association rules: Apriori implementation. In Proceedings of the 15th Conference on Computational Statistics (Compstat 2002, Berlin, Germany) (pp. 395-400). Physica-Verlag HD, 2002.
- [84] C. Borgelt. Efficient implementations of apriori and eclat. In FIMI'03: Proceedings of the IEEE ICDM Workshop on Frequent Itemset Mining Implementations (FIMI 2003, Melbourne, FL, USA, 2003).
- [85] C. Borgelt. Recursion Pruning for the Apriori Algorithm. In FIMI'04: Proceedings of the 2nd IEEE ICDM Workshop of Frequent Item Set Mining Implementations (FIMI 2004, Brighton, UK), 2004.
- [86] D. Hansen, B. Shneiderman, M.A. Smith. *Analyzing social media networks with NodeXL: Insights from a connected world*. Morgan Kaufmann, 2010.



Xu Chi Dr. Xu Chi received his Ph.D. and Bachelor (Honors) in Electrical and Electronic School from Nanyang Technological University, Singapore, in 2010 and 2003 respectively. He was a researcher in Positioning and Wireless Technology Center in Nanyang Technological University, working on RFID ranging and positioning using ultra wideband (UWB) signal. Currently, he is a research scientist in the Planning and Operations Management Group of Singapore Institute of Manufacturing Technology (SIMTech). His research interests include information management for track and trace system and unstructured data mining.



Allan N. Zhang Dr Allan N. Zhang is a Senior Scientist with Singapore Institute of Manufacturing Technology, A*STAR, Singapore. He has more than 20 years experience in knowledge-based systems and enterprise information systems development. His research interests include knowledge management, data mining, machine learning, artificial intelligence, computer security, software engineering, software development methodology and standard, and enterprise information systems. He and his team are currently working toward research in manufacturing system analyses including data mining, supply chain information management, supply chain risk management using Complex Systems approach, and urban last mile logistics.



Gürdal Ertek Dr. Gürdal Ertek is an Assistant Professor at Rochester Institute of Technology - Dubai. Earlier, he was with Sabanci University, Istanbul, Turkey and was a Visiting Scientist at Singapore Institute of Manufacturing Technology (A*Star SIMTECH). He received his B.S. from Industrial Engineering Department of Bogazici University, Istanbul, Turkey, in 1994, and his Ph.D. from School of Industrial and Systems Engineering at Georgia Institute of Technology, Atlanta, GA, in 2001. He has been awarded with Bogazici University Alumni Scholarship, Haci Omer Sabanci Scholarship, and Fulbright Scholarship throughout his education. His research areas include knowledge-based systems, warehousing and material handling, and data visualization and mining. Dr. Ertek has served as a reviewer for 50+ R&D projects submitted to TUBITAK (Turkish National Science Foundation), mostly on the topics of information technology and data analytics.