

# Temporal Knowledge Discovery in Big BAS Data for Building Energy

## Management

Cheng Fan<sup>a</sup>, Fu Xiao<sup>a,\*</sup>, Henrik Madsen<sup>b</sup>, Dan Wang<sup>c</sup>

<sup>a</sup>Department of Building Services Engineering

<sup>c</sup>Department of Computing

The Hong Kong Polytechnic University, Hong Kong

<sup>b</sup>Department of Applied Mathematics and Computer Science

Technical University of Denmark

\* E-mail: [linda.xiao@polyu.edu.hk](mailto:linda.xiao@polyu.edu.hk); Tel.: (852) 27664194

## Abstract

With the advances of information technologies, today's Building Automation Systems (BASs) are capable of managing building operational performance in an efficient and convenient way. Meanwhile, the amount of real-time monitoring and control data in BASs grows continually in the building lifecycle, which stimulates an intense demand for powerful big data analysis tools in BASs. Existing big data analytics adopted in the building automation industry focus on mining cross-sectional relationships, whereas the temporal relationships, i.e., the relationships over time, are usually overlooked. However, building operations are typically dynamic and BAS data are essentially multivariate time series data. This paper presents a time series data mining methodology for temporal knowledge discovery in big BAS data. A number of time

series data mining techniques are explored and carefully assembled, including the Symbolic Aggregate approXimation (SAX), motif discovery, and temporal association rule mining. This study also develops two methods for the efficient post-processing of knowledge discovered. The methodology has been applied to analyze the BAS data retrieved from a real building. The temporal knowledge discovered is valuable to identify dynamics, patterns and anomalies in building operations, derive temporal association rules within and between subsystems, assess building system performance and spot opportunities in energy conservation.

**Keywords:** Temporal knowledge discovery; Time series data mining; Big data; Building automation system; Building energy management

## **1. Introduction**

The building sector is evolving to be the greatest energy consumer around the world, accounting for 40% of the global energy use and one third of the global greenhouse gas emissions [1, 2]. As a result, building energy efficiency has become one of the top concerns of a sustainable society and attracted increasing research and development efforts in recent years. Thanks to the advances of information, computing and control technologies, Building Automation System (BAS) provides a valuable network-based digital platform for automatically managing complex building systems, including heating, air conditioning, ventilation, lighting, vertical transportation, fire safety and security systems. It is estimated that the potential energy savings from the adoption of

1 advanced building automation technologies might reach 22% by 2028 for the  
2 European building sector [3]. Besides fulfilling the online monitoring and control  
3 functions, BASs record thousands of real-time measurements and control signals, and  
4 the amount of data keeps growing in the building life-cycle. Most of the existing  
5 building management strategies are developed based on domain expertise or small  
6 subsets of the BAS data. There are increasing interests in systematically analyzing the  
7 big BAS data and discover knowledge for improving building performance.

8 Data mining (DM) is a promising technology, which can effectively discover  
9 interesting and potentially useful knowledge from big data. Some efforts have been  
10 made to investigate the potentials of DM in the building field. DM techniques have  
11 been adopted at the building design, construction, and operation stages [4]. The  
12 building operation stage draws particular attention, as it accounts for 80-90% of the  
13 total building green gas emission and is directly linked to occupant comforts and the  
14 realization of building functionality [5]. In general, DM techniques can discover two  
15 types of knowledge, i.e., predictive and descriptive. Predictive DM is often used to  
16 capture the complex and nonlinear relationships between inputs and outputs. It has  
17 been applied at the building operation stage to the prediction of building energy  
18 consumption [6, 7, 8], thermal load [9, 10], indoor environment [11, 12], and system  
19 performance indices [13, 14, 15]. Descriptive DM is used to discover the associations,  
20 correlations, and intrinsic data structure in big data. Compared to predictive DM,  
21 descriptive DM is more flexible in applications, as it does not involve a training  
22 process and the knowledge discovery process is not guided by pre-defined targets.

1 Descriptive DM has been mainly applied at the building operation stage to fault  
2 detection and diagnostics [4, 16, 17, 18]. Popular techniques include association rule  
3 mining, clustering analysis, and anomaly detection.

4 Despite of the encouraging research outcomes, previous studies of analyzing big BAS  
5 data usually considered the BAS data as cross-sectional data, and hence the  
6 knowledge discovered mainly includes the concurrent relationships among different  
7 variables, such as relationships between the power consumptions of the primary air  
8 units and lifts in a building [18]. BAS data are usually recorded in a two-dimensional  
9 matrix, with each column representing a measured variable (such as temperature, flow  
10 rate, power and control signal), and each row representing an observation/sample at a  
11 specific time instant. Each observation is a vector of many measurements and control  
12 signals. The time interval between two consecutive observations is usually fixed and  
13 may vary from several seconds to tens of minutes. The first two columns usually store  
14 the date and the time. Considering that BAS data are in essence multivariate time  
15 series data, the cross-sectional knowledge discovered may not be able to fully capture  
16 the relationships over time. Building operations are typically dynamic due to the  
17 changes in indoor and outdoor operating conditions, such as the outdoor climate  
18 conditions, indoor occupant number and utilization of indoor electric appliances.  
19 Meanwhile, the changes hardly occur simultaneously which results that the dynamics  
20 in building operations are very complicated. For instance, the indoor temperature is  
21 influenced by the outdoor air temperature. However, when the infiltration is not  
22 significant, these two temperatures rarely change simultaneously due to building

1 thermal mass. Time lags between them often bring challenges to the sequence control  
2 of chiller plants. The dynamics are usually complicated and have great influences on  
3 control performance, interactions among building components and integrations  
4 between buildings and communities (e.g., electricity power grid) [19]. In practice, it is  
5 desired to discover such temporal knowledge hidden in BAS data. Advanced tools  
6 and methods for temporal knowledge discovery should be developed for this purpose.  
7 Conventional time series analytics, such as the autoregressive moving average models  
8 (ARMA), are mainly used for solving predictive tasks in the field of building  
9 management, including the prediction of building electricity consumption [20, 21],  
10 building thermal load [22, 23] and indoor environment [24, 25]. In recent years,  
11 various approaches have been developed to mine temporal knowledge in different  
12 formats, such as events, clusters, motifs and temporal association rules [26, 27, 28,  
13 29, 30]. However, only limited studies have been performed to explore their potential  
14 in analyzing BAS data. The complex event processing (CEP), which is a method well  
15 suited for the processing of information flows, has been adopted to utilize time series  
16 data in building operations [26, 27, 28]. Renner et al. applied CEP to correlate the  
17 sensor data and provide real-time reactions to building management [26]. Wen et al.  
18 developed a CEP-based method to derive knowledge from building energy data and  
19 applied it for building controls [27]. In these studies, domain expertise plays an  
20 important role in defining events and rules. Time series data mining enables an  
21 approach to discover interesting and previously unknown temporal knowledge.  
22 Patnaik et al. adopted the motif discovery technique to mine chiller operation data in

1 data centers [31]. Motifs (i.e., frequent sequential patterns) were successfully  
2 discovered to identify energy-efficient operating patterns. Miller, Nagy and Schlueter  
3 used a similar method to analyze building energy consumption data [32]. Energy  
4 consumption motifs were extracted for building performance characterization.  
5 Discords, or infrequent sequential patterns, were identified and used for fault  
6 detection. Their work demonstrated the encouraging potentials of time series data  
7 mining in the knowledge discovery of BAS data for managing building operations.  
8 Currently, the potential and applicability of various time series data mining  
9 techniques in mining big BAS data are still uncertain considering unique  
10 characteristics of BAS data, such as low quality, nonlinearity, multiple scales or units,  
11 and multicollinearity. A generic and systematic methodology for discovering  
12 temporal knowledge in big BAS data is needed for developing applicable tools in  
13 BAS.

14 This study proposes a generic methodology for mining temporal knowledge hidden in  
15 big BAS data and demonstrates its applications in real cases. We first briefly  
16 introduce time series data mining techniques, and specifically highlight the  
17 differences between cross-sectional data mining and time series data mining. Then,  
18 the generic methodology is presented, which consists of data preprocessing, data  
19 partitioning, temporal knowledge discovery, and post-mining. Two methods are  
20 developed to improve the efficiency in post-mining. The methodology has been  
21 applied to analyze the BAS data retrieved from the tallest building in Hong Kong.  
22 Valuable temporal knowledge has been discovered for building operations and

1 performance management.

2

## 3 **2. Description of Research Methodology**

4 Based on a comprehensive exploration of advanced DM techniques, in-depth analysis  
5 of BAS data characteristics as well as specific considerations for practical  
6 applications, we have developed a generic framework for knowledge discovery in  
7 BAS data [4]. The framework consists of four major phases, i.e., data preprocessing,  
8 data partitioning, knowledge discovery and post-mining. Each phase was specifically  
9 designed considering the BAS data quality and structure, data format requirement of  
10 DM techniques, interpretation and selection of knowledge discovered, and application  
11 of the knowledge to building performance assessment, diagnosis and optimization. It  
12 is a generic framework proposed for analyzing big BAS data using DM techniques.  
13 The methodology presented in this study is also developed within this framework, as  
14 shown in Figure 1, and further enriches the framework by integrating time series data  
15 mining techniques for temporal knowledge discovery. Three tasks are performed in  
16 the first phase, including data cleaning, period estimation and data transformation.  
17 Phase 2 adopts the evidence accumulation clustering to partition the SAX  
18 subsequences. Phase 3 adopts two techniques, i.e., motif discovery and temporal  
19 association rule mining, to discover two different types of knowledge. Two  
20 post-mining methods are developed in Phase 4 to improve the efficiency and  
21 effectiveness of handling the large amount of knowledge discovered in Phase 3. The  
22 details of each phase are introduced in the following subsections.

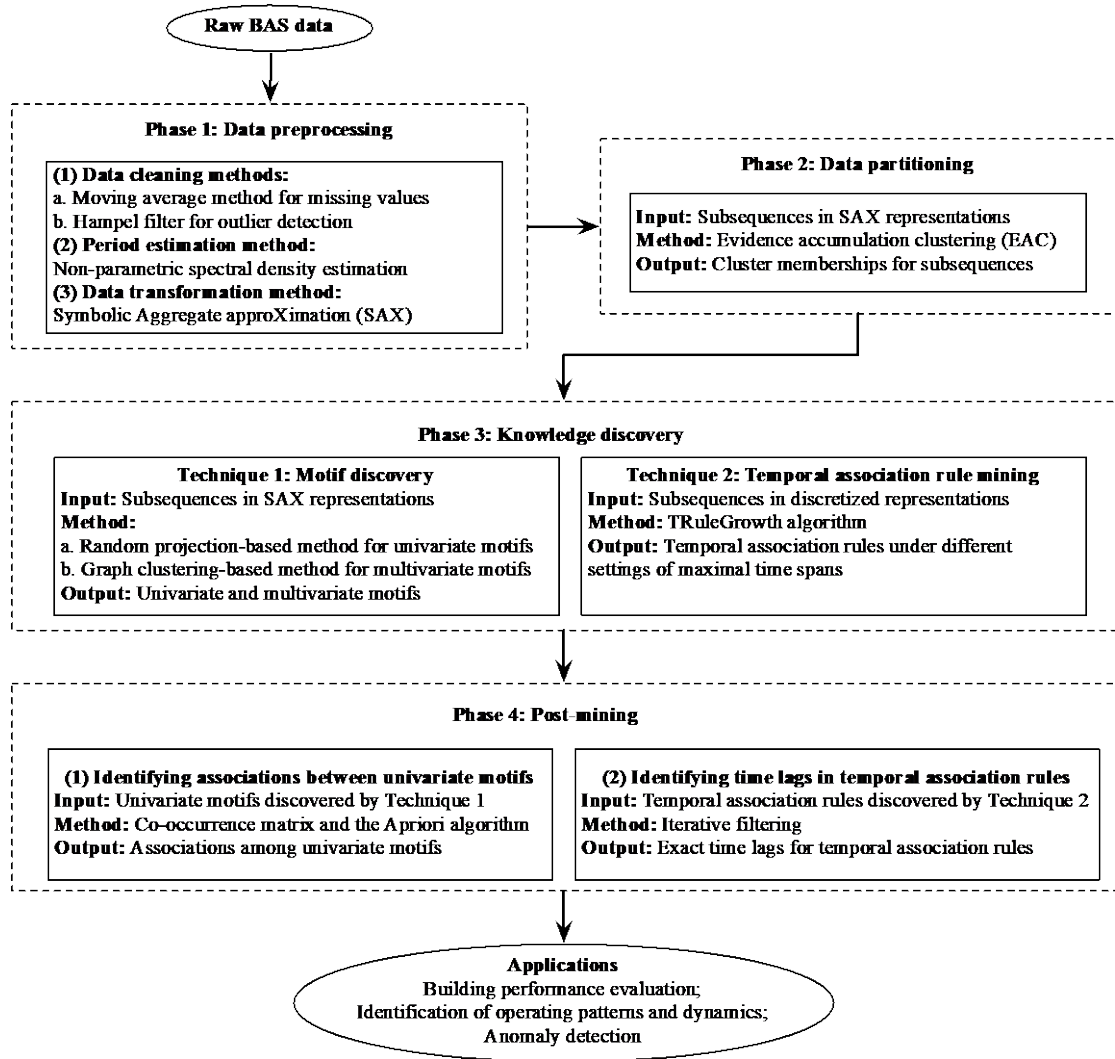


Figure 1 Outline of the research methodology

## 2.1 Data preprocessing

Data preprocessing fulfills three tasks, i.e., data cleaning, period estimation, and data transformation, with the aims to enhance the data quality, explore the intrinsic characteristics in BAS time series data, and prepare the raw data with suitable format for data mining.

### 2.1.1 Data cleaning



1 Data cleaning aims to improve BAS data quality by filling missing values and  
2 detecting outliers in raw BAS time series data. Missing values widely exist in BAS  
3 data due to sensor malfunctions or signal transmission problems in BAS network.  
4 Popular methods to impute missing data include moving window, random imputation,  
5 and inference-based methods [33]. Moving window-based methods are easy to  
6 implement and have a fairly good performance when the duration of missing values is  
7 not very long. Otherwise, it is recommended to use the inference-based methods or  
8 simply exclude the observations with long missing values from analysis. Outliers in a  
9 time series are observations that are highly unlikely to occur based on the variation  
10 seen in the rest of the time series. They can be classified into two types, i.e., points as  
11 outliers and subsequence as outliers [34]. It is recommended that the data  
12 preprocessing phase only handles the first type of outliers in the raw time series data,  
13 as the identification of the second type of outliers may overlap with mining discords  
14 (i.e., infrequent sequential patterns) in the later process. The outlier detection methods  
15 can be grouped into three categories, i.e., prediction-based, profile-based, and  
16 deviant-based methods [34]. The prediction-based methods detect outliers by  
17 comparing the actual measurements with their expected or predicted values from  
18 statistical analysis or machine learning algorithms. The profile-based methods use  
19 historical data to construct a normal profile, which is usually presented in the form of  
20 expected means and confidence intervals at different time. Each observation is  
21 compared with the normal profile to decide whether it is an outlier or not. The  
22 deviant-based methods identify outliers from a perspective of information theory. An

1 observation is an outlier if removing it from the time series leads to a much more  
2 succinct representation of the original time series [34].

3 In this study, the moving window-based method is used to impute the missing values  
4 with a short duration, i.e., less than 2-hour. Any missing values with a longer time  
5 duration are excluded from analysis. This study adopts the Hampel filter to identify  
6 outliers. It is a nonlinear filter which shows high effectiveness in processing time  
7 series data [35]. For each observation, the Hampel filter calculates the median and the  
8 median absolute deviation (MAD) considering a moving window size of  $2k+1$ .  $k$  is  
9 the number of observations before and after the observation concerned. A  
10 parameter  $\theta$ , which usually ranges from 0 to 5, is predefined to generate thresholds  
11 for outlierness evaluation, i.e.,  $Median \pm \theta \times MAD$ . Any observation falls beyond  
12 the range of the thresholds is identified as an outlier and is replaced by the median.  
13 The smaller the parameter, the more aggressive the detection algorithm is and more  
14 observations will be identified as outliers. This study sets  $\theta$  as 3, which is in  
15 accordance with the Ron Pearson's 3-sigma rule.

16

### 17 **2.1.2 Period estimation**

18 This step is specifically developed for time series data, considering that long time  
19 series data usually exhibit periodicity, and consequently motifs and association rules  
20 periodically repeat. Finding the period in the time series data and then segment those  
21 data into short subsequences can considerably reduce the mining load. It is a common  
22 practice in time series data mining, particularly in handling very long time series data

like BAS data. The repeating daily working schedule of building users (e.g. office hours and non-office hours) results that the operating schedules of major systems and equipment (such as air conditioning, lighting and lift systems) usually repeat daily. Obviously, the BAS data exhibit daily periodicity. In view of this, Miller, et al. segmented the time series of building energy consumption data into daily sequences in their study [32]. This study attempts to adopt a data-driven approach to estimating the intrinsic periods embedded in BAS data. There are two purposes for doing this: firstly to minimize the dependence on domain knowledge in the knowledge discovery process; secondly to maximize the possibility of discovering new knowledge, or new periods in BAS data in our case. Periods in time series data can be detected using the spectral density estimation methods, which can be either parametric or non-parametric. The parametric methods first model the time series using time series modeling techniques, such as autoregressive and moving average (ARMA). The spectral density is then estimated based on the model parameters. By contrast, the non-parametric methods estimate the spectral density by taking the Fourier transformation of the autocorrelation function. Considering that the building data usually present diurnal, weekly and annual seasonality, the resulting parametric models can be very complex. Therefore, this study applies the non-parametric method to period estimation.

### **2.1.3 Data transformation**

Data transformation prepares the time series data with suitable formats to meet the

1 following two needs. Firstly, different mining techniques require different data  
2 formats (e.g., numerical or categorical) and BAS data exhibit diversity in units, scales,  
3 and data types. Secondly, the computation load is a big concern due to the huge  
4 volume of big data, which can be alleviated by effectively reducing the volume of the  
5 data without losing valuable information embedded in the data. In this study, the  
6 symbolic approximation aggregate (SAX) method is proposed to transform the  
7 original time series BAS data into meaningful symbols [32, 36]. The SAX method  
8 transforms a numeric time series into a symbol stream and the length of the symbol  
9 stream is much shorter than the original time series. It can therefore reduce the data  
10 size.

11 To perform SAX, a univariate time series of length  $n$  is firstly standardized to have a  
12 zero mean and a standard deviation of 1 and then segmented into  $m$  subsequences  
13 with a window size of  $q$ . One of the typical methods to segment the time series is  
14 based on the period detected in the previous step. For example, if the period estimated  
15 is 24 hours, one day BAS data will form one subsequence. Two parameters need to be  
16 defined to perform SAX, i.e., the word size  $W$  and the alphabet size  $A$ . A set of  
17 breakpoints (e.g.,  $\beta_1, \beta_2, \dots, \beta_{A-1}$ ) are determined in such a manner that the area  
18 under the  $N(0,1)$  Gaussian curve from  $\beta_i$  to  $\beta_{i+1}$  is  $\frac{1}{A}$ . Each interval will be  
19 assigned with an alphabet (e.g.,  $a$ ,  $b$ , and  $c$ ) and the number of alphabets used is the  
20 alphabet size,  $A$ . Given the word size ( $W$ ), each subsequence in the window size of  $q$   
21 can be divided into  $W$  equal sections, and the means of each sections are calculated.  
22 According to which interval (i.e.,  $\beta_i$  to  $\beta_{i+1}$ ) the mean lies within, the corresponding

1 alphabet is assigned to the section. In this way, each subsequence can be represented  
2 by a SAX word which consists of  $W$  alphabets. For example, *abca*, *aabc*, *bcca* are  
3 SAX words given  $W=4$  and  $A=3$ . In these SAX words, the alphabet size ( $A$ ) is 3, so  
4 three alphabets (i.e.,  $a$ ,  $b$  and  $c$ ) are used; the word size ( $W$ ) is 4, so each SAX word  
5 consists of four alphabets. The original time series is transformed into a string of  
6 alphabets. The larger the alphabet size ( $A$ ) and the word size ( $W$ ), the more detailed  
7 information retained in the symbolic stream. However, the reduction of computation  
8 load becomes less. Therefore, there is a trade-off, which will be discussed in the later  
9 case studies.

10 The distance between two SAX representations are calculated as  $\sqrt{\frac{q}{w}} \times$   
11  $\sqrt{\sum_{i=1}^w dist(S_i, B_i)^2}$ , where  $S$  and  $B$  are two SAX representations, and  $dist()$  is the  
12 distance function for SAX symbols. Table 1 presents an example of distance matrix  
13 between symbols considering an alphabet size of 4. The value in  $cell(x,y)$  is calculated  
14 using Equation 1. A dissimilarity matrix considering different SAX representations  
15 can be computed accordingly. More details can be found in [36].

16 Besides SAX, difference-based and dictionary-based methods are also capable of  
17 transforming time series into symbols [37, 38, 39]. The difference-based method  
18 transformed the raw time series into symbols based on their first- or higher-order  
19 differences. It can be used when the changes between successive time steps are more  
20 important than the absolute values [39]. The dictionary-based methods transform the  
21 time series into symbols by matching the raw data with predefined patterns in a  
22 dictionary. For instance, in the studies performed by Kwac et al. [37] and Gulbinas et

al. [38], clustering analysis was applied to generate the representative patterns of daily power consumption, based on which a dictionary was built for symbolization. SAX is selected in this study considering the following two aspects. Firstly, SAX is straightforward to use, as it requires little domain expertise and preprocessing. Secondly, SAX contains an intrinsic distance measure, which provides extra value in the subsequent knowledge discovery [36], as shown in the later part.

$$\text{Equation 1: } \quad cell(x, y) = \begin{cases} 0 & \text{if } |x - y| \leq 1 \\ \beta_{\max(x,y)-1} - \beta_{\min(x,y)} & \text{otherwise} \end{cases}$$

Table 1: An example distance matrix for SAX symbols

Distance	$a$	$b$	$c$	$d$
$a$	0	0	0.67	1.34
$b$	0	0	0	0.67
$c$	0.67	0	0	0
$d$	1.34	0.67	0	0

## 2.2 Data partitioning

Due to the changing operating conditions and complicated system dynamics and interactions, the big BAS data usually scatter in a high-dimensional space. To enhance the reliability and sensitivity of the mining results, data partitioning is carried out to divide the data into several groups or clusters, with the aim of maximizing the intra-group similarities while minimizing the inter-group similarities. Knowledge discovery are then performed on each group separately. Clustering analysis is a suitable DM technique to perform this task. Despite of the large number of clustering

1 algorithms being available, no single algorithm is able to identify all kinds of cluster  
2 shapes and data structures in practice [40]. It is usually very difficult to find out the  
3 optimal clustering algorithm and the settings of its parameters. Some methods have  
4 been developed to facilitate the decision-makings, based on either internal (e.g., Dunn  
5 index and Davies-Bouldin index) or external validation indices (e.g., purity and  
6 mutual information). However, no validation method can impartially evaluate the  
7 results of any clustering algorithm [41]. A common practice is to try out a large  
8 number of algorithms with different parameters in order to obtain desired the  
9 clustering results. The process can be computationally expensive and  
10 time-consuming.

11 Ensemble learning is capable of enhancing the clustering performance by combining a  
12 number of base learners, whose individual performance may be poor [40, 41]. The  
13 evidence accumulation clustering (EAC) is a method designed to apply ensemble  
14 learning on clustering analysis [34]. One advantage of the EAC over other  
15 conventional clustering methods is that it has the ability to discover clusters with  
16 various sizes and shapes. In addition, the method can automatically determine the  
17 optimal cluster number, which provides great flexibility in analyzing data with  
18 unknown characteristics. The partition around medoids (PAM) is selected as the base  
19 algorithm for EAC. PAM shares a similar partitioning mechanism as the popular  
20  $k$ -means algorithm. Compared to the  $k$ -means, PAM is more robust to outliers and  
21 noises and can take a dissimilarity matrix as inputs. Therefore, PAM is more  
22 compatible with time series data in SAX representations.

1 Three parameters needs to be defined to perform EAC, i.e., the total iteration number  
2  $E$ , the lower and upper limits of the cluster number  $K_{lower}$  and  $K_{upper}$ .  $E$  sets of  
3 clustering results are generated by PAM with different cluster numbers (i.e., randomly  
4 sampled from  $K_{lower}$  to  $K_{upper}$  in each iteration) and the dimension of input data. These  
5  $E$  sets of clustering results are then transformed into a co-occurrence matrix.  
6 Assuming that the data contains  $n$  observations, the co-occurrence matrix  $C$  has a  
7 dimension of  $n \times n$ . The value of  $C_{ij}$  is the number of times when observations  $i$  and  
8  $j$  are grouped in the same cluster divided by the total iteration number  $E$ . The final  
9 clustering result is obtained by using hierarchical agglomerative method to cluster the  
10 co-occurrence matrix. More details can be found in [40].

11

## 12 **2.3 Temporal Knowledge Discovery**

13 After the data are preprocessed and partitioned, appropriate DM techniques will be  
14 applied for knowledge discovery. The typical descriptive knowledge types in time  
15 series data include motifs, discords and temporal association rules [29].

16

### 17 **2.3.1 Motif discovery**

18 Motif, or frequent sequential pattern, is a typical knowledge type which can be  
19 discovered in time series data. Motifs are valuable to temporal association rule  
20 mining, discord (i.e. infrequent sequential pattern) detection, and time series  
21 classification [42].

22 Motif discovery has been mainly applied to analyze univariate time series in previous



1 studies. Conventional motif discovery methods are based on exhaustive search, which  
 2 results that the computational costs increase dramatically for long time series and is  
 3 therefore not applicable to big data. In view of this, a more efficient algorithm, which  
 4 is based on random projection and compatible with SAX representations [42], is  
 5 selected to discover univariate motifs. Assuming that the time series has a length of  $n$   
 6 and the sliding window size is  $q$ , a matrix containing all the subsequences (denoted as  
 7  $M_1$ ) can be constructed and has a dimension of  $(n - q + 1) \times q$ . Each subsequence is  
 8 transformed into a SAX representation. Assuming the word size is  $W$ , the new matrix  
 9 containing the SAX representations (denoted as  $M_2$ ) has a dimension of  $(n - q +$   
 10  $1) \times W$ . Random projection is performed by randomly picking  $s$  columns from  $M_2$ ,  
 11 where  $s$  ranges from 1 to  $W-1$ . A collision matrix, which has a dimension of  $(n - q +$   
 12  $1) \times (n - q + 1)$ , is constructed to record the times of being identical for two  
 13 subsequences after a number of random projections. A tentative univariate motif is  
 14 identified if the two subsequences result in a high value in the collision matrix.  
 15 Potential members of this tentative univariate motif can then be identified by  
 16 calculating the Euclidean distance in the original numeric representations.  
 17 Several methods have been developed to identify motifs in multivariate time series  
 18 data, such as PCA-based and density estimation-based methods [43, 44]. Those  
 19 methods can successfully identify synchronous multivariate motifs. However, their  
 20 practical value in analyzing real-world data is limited, as the motifs in multivariate  
 21 time series data do not necessarily start at the same time and their duration may vary  
 22 as well. We can see a lot of such examples in building operations. For example, when

the air conditioner or chiller is turned on, the indoor temperature will not change immediately due to the thermal mass. The sudden increase of the lift power consumption in the morning peak hour does not correspond to a large increase in the chiller power consumption due to the pre-cooling strategy. In this study, multivariate motif discovery algorithm proposed in [45] is adopted. The main advantage is that, firstly, both synchronous and non-synchronous multivariate motifs can be discovered, and secondly, the multivariate motifs identified may consist of all univariate motifs or any subset of the univariate motifs. The method first performs univariate motif discovery on the time series of each variable. A graph clustering approach is then applied to identify multivariate motifs. A directed coincidence graph  $G$  is constructed. Each motif  $r_i$  is represented by a vertex  $v_i$ .  $e_{i,j}$  represents the edge connecting the vertex  $v_i$  and  $v_j$ . The weight of  $e_{i,j}$  is denoted as  $w_{i,j}$  and calculated as  $\text{coincident}(r_i, r_j)/\text{size}_i$ , where  $\text{coincident}(r_i, r_j)$  is the total number of times that a temporal overlap is found between  $r_i$  and  $r_j$  and the  $\text{size}_i$  is the number of occurrence of  $r_i$ . A parameter,  $\alpha$ , ranging from 0 to 1, is user-specified as the minimum correlation between univariate motifs based on which a multivariate motif could be constructed.

### 2.3.2 Temporal association rule mining

The difference between association rule mining (ARM) and temporal association rule mining (TARM) lies in whether the temporal information is contained in the rule or not. ARM was mainly used to discover cross-sectional associations, where the temporal information is neglected. The typical format of ARM is  $A \rightarrow B$ , where  $A \cap$

1  $B = \emptyset$ . It states that if  $A$  happens,  $B$  will also happen. An association rule is derived if  
 2 both the rule support and confidence exceed the user-defined thresholds. The support  
 3 of a rule is the fraction between the number of times when both the antecedent and  
 4 consequent take place and the total number of records. The confidence of a rule is the  
 5 conditional probability of the consequent given the antecedent. The interestingness of  
 6 the association rules can be evaluated using the *lift*, which is the ratio between the rule  
 7 confidence and the support of consequent. It measures the dependency and correlation  
 8 between the antecedent and the consequent of a rule. Potentially useful rules usually  
 9 have a *lift* larger than 1, indicating that the occurrence of the antecedent positively  
 10 influences the occurrence of consequent.

11 Temporal association rule mining (TARM) is of particular interest in mining BAS  
 12 data because of the complicated dynamics in building operations. TARM, or  
 13 sequential rule mining, discovers associations among variables while providing an  
 14 insight into the temporal dependency. The general format of temporal association  
 15 rules is also  $A \rightarrow B$ , where  $A \cap B = \emptyset$ . However, the temporal dependency is  
 16 contained, indicating that  $B$  will take place after  $A$ . Various algorithms have been  
 17 developed for deriving temporal association rules, such as the SPADE and CMRules  
 18 [46, 47]. In engineering practice, temporal rules that are valid within a limited time  
 19 span are of special interest. The format of such temporal rules is  $A \xrightarrow{t} B$ , which means  
 20 that  $B$  will occur within  $t$  time units after the occurrence of  $A$ . Therefore, the  
 21 TRuleGrowth algorithm, which can derive temporal association rules under the  
 22 constraint of maximum time span [48], is selected in this study. To perform this

1 algorithm, three parameters need to be defined, i.e., the minimum support, minimum  
2 confidence, and the maximum time span. The other advantage of the TRuleGrowth  
3 algorithm is that it can greatly reduce the number of rules generated by controlling the  
4 maximum time span. Consequently, the post-mining phase consumes much less time.

5

## 6 **2.4 Post-mining**

7 The post-mining phase aims to build a bridge between knowledge discovered in Phase  
8 3 and practical applications, such as building performance assessment, fault diagnosis  
9 and optimization. It usually needs domain expertise to select, interpret and apply the  
10 knowledge discovered [4, 18, 32]. The process can be very time-consuming, due to  
11 the large amount of knowledge discovered and the diversity of knowledge  
12 representations (e.g., rules, clusters, decision trees). Application of the motifs and  
13 temporal association rules is straightforward. They can be used as the references for  
14 normal operations and anomalies can be detected if building operation patterns are  
15 different from those frequent patterns or violate the association rules. In this study,  
16 two methods are specifically developed to enhance the efficiency in post-mining and  
17 maximize the practical values of temporal knowledge discovered.

18

### 19 **2.4.1 Identify associations between univariate motifs**

20 Building operations involves multiple separate and interactive subsystems, such as air  
21 conditioning, mechanical ventilation, lift, lighting and security systems. Univariate  
22 motifs usually represent the frequent sequential operation patterns of each system. It

1 is reasonable to link the associations among univariate motifs with the interactions  
 2 among subsystems. Multivariate motifs can provide general information on which  
 3 univariate motifs frequently occur together. However, they hardly quantify the  
 4 relationships among univariate motifs and this limits their practical value. For  
 5 example, a multivariate motif cannot answer, if one univariate motif occur, whether  
 6 the other univariate motifs in it will occur or not with certain probability. In this  
 7 study, a post-mining method is designed to explore the associations among univariate  
 8 motifs which can directly answer this question. This method is an extension of  
 9 association rule mining. Given a multivariate time series data, the univariate motif  
 10 discovery algorithm is applied to each univariate time series separately to find  
 11 univariate motifs. These univariate motifs are then labeled as  $m_1, m_2, \dots, m_L$ , where  $L$   
 12 is the total number of univariate motifs discovered. Afterward, a co-occurrence matrix  
 13 is constructed. The matrix has  $L$  columns. The values of each row are either 1 or 0,  
 14 indicating whether an occurrence of a univariate motif is observed or not. Once the  
 15 matrix is constructed, the Apriori algorithm is used to discover associations between  
 16 univariate motifs. Two parameters, i.e., the minimum thresholds for support and  
 17 confidence, are defined for rule induction. Three statistics, including the support,  
 18 confidence and lift, can be generated with each association rule to facilitate decision  
 19 making.

20 An example for construction of a co-occurrence matrix is given here. Figure 2  
 21 illustrates five univariate motifs (i.e.,  $m_1$  to  $m_5$ ) in the sequences of three variables,  $A$ ,  
 22  $B$  and  $C$ . The motifs in  $A$  and  $B$ ,  $m_1$  to  $m_4$ , have a time duration of 10 while  $m_5$  in  $C$

1 has a time duration of 8. The co-occurrence matrix is constructed as shown in Table 2.  
 2 The numbers (0 or 1) in each row show the occurrence of the corresponding motifs.  
 3 For example, the second row shows that only motif 5 occurs during the time period  
 4 between 18 to 25; the first and third rows show that motifs 1 and 3 occur together  
 5 twice; the fifth row shows that motifs 2, 4, 5 occur together for once; the sixth row  
 6 shows that motifs 2, 3 and 5 occur together for once. It should be noted that, although  
 7 the occurrence of the motifs are related to certain time period, the exact time is not  
 8 considered in constructing the matrix. The frequency of the co-occurrence of multiple  
 9 univariate motifs is of interest.

10 The construction of the co-occurrence matrix can be conveniently implemented by  
 11 programming with the information of starting and ending time instants of all  
 12 univariate motifs. Once the co-occurrence matrix is ready, the Apriori algorithm is  
 13 adopted to mine the associations. Setting the minimum thresholds of support and  
 14 confidence as 0.3 and 0.8 respectively, two rules are derived, i.e.,  $m_1 \rightarrow m_3$  and  
 15  $m_2 \rightarrow m_5$ . Both rules have a support of 0.4 and a confidence of 1. It means that when  
 16 motif 1 occurs, the probability of the occurrence of motif 3 is very high.

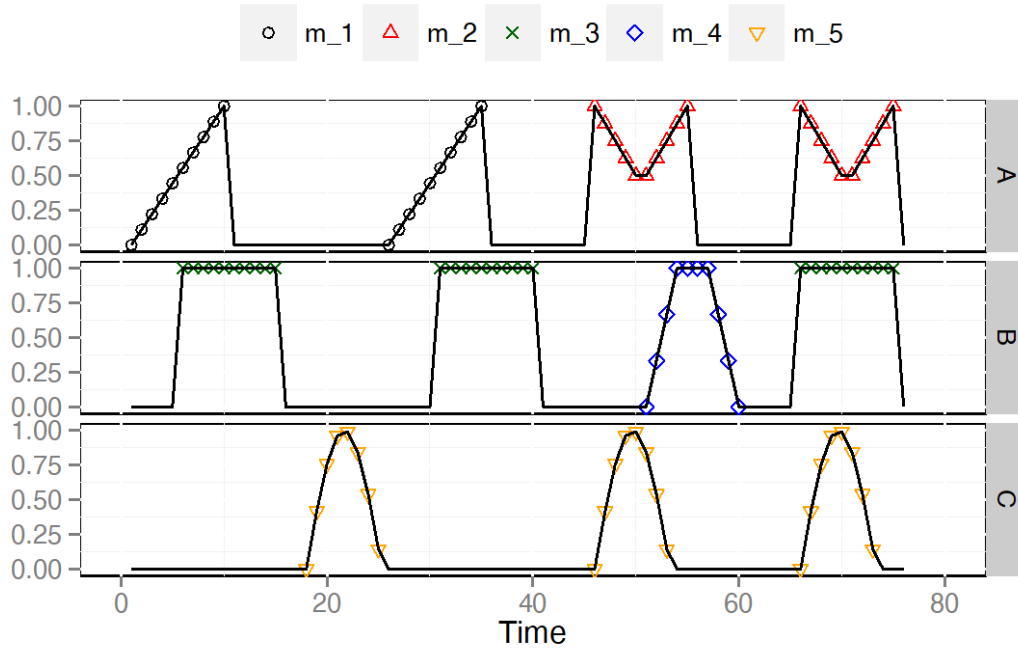


Figure 2: An example of univariate motifs discovered in three dimensions

Table 2: An example of co-occurrence matrix for mining association rules between univariate motifs

$m_1$	$m_2$	$m_3$	$m_4$	$m_5$
1	0	1	0	0
0	0	0	0	1
1	0	1	0	0
0	1	0	1	1
0	1	1	0	1

#### 2.4.2 Identify time lags in temporal association rules

As introduced in Section 2.3.2, the TRuleGrowth algorithm is adopted to discover the temporal association rules under the constraint of a maximum time span. One limitation is that no information is available about the exact time lag, which is the time interval between the antecedent and the consequent. This type of information is valuable for establishing reliable control and performance optimization in building

1 operations. An iterative filtering method is developed to identify the time lag. The  
2 method iteratively runs the TRuleGrowth algorithm by changing the maximum time  
3 span from 1 to  $T$  and the temporal association rules generated at each iteration and the  
4 corresponding time lag are stored in the rule sets. The time lag in a temporal  
5 association rule can be discovered by matching the rule with the rule sets.

6

### 7 **3. Mining Real BAS Data**

#### 8 **3.1 BAS data description**

9 The methodology is applied to analyze the BAS data retrieved from the tallest  
10 building in Hong Kong, i.e., the International Commerce Center (ICC). ICC is a  
11 high-rise commercial building with a height around 490m and a total floor area of  
12 321,000m<sup>2</sup>. It contains a 4-storey basement for parking, a 6-storey block for shopping  
13 and exhibitions, a 65-storey office tower, an observations deck and a 17-storey hotel.  
14 A complex BAS has been installed to monitor and control the building operations.  
15 Energy efficiency is one of the major concerns of ICC. The whole building power  
16 consumption can be broken down into five parts for different services systems, i.e.,  
17 the heating, ventilation, and air-conditioning (HVAC) system, normal power and  
18 lighting (NLTG), essential power and lighting (ELTG), vertical transportation system  
19 (VTS), and plumbing and drainage system (PD). The HVAC system in ICC consists  
20 of six subsystems, i.e., chillers, cooling towers, water pumps, primary air-handling  
21 units (PAU), air-handling units (AHU), and mechanical ventilation (MV). Besides  
22 power consumption data, there are a large number of measurements of temperature,



flow rate, pressure, etc., of the water and the process air in the system as well as various status and control signals. There are approximately 950 measurements in the BAS system, which are sampled every 1 or 15 minutes. The size of annual BAS data in ICC is around 30 gigabytes. Annual operation data in 2014 with a sampling interval of 15-minute are retrieved for analysis in this study. The data consist of 34,950 observations and 78 variables, including the date and time, building cooling load as well as power consumptions of various subsystems. Approximately 0.52% of the data contain missing values and the maximum lasting period is 45-minute. The moving average method with a window size of 8 (i.e., corresponding to 2-hour) is adopted to fill the missing values. The Hampel filter method is applied to detect point-wise outliers and parameters  $k$  and  $\theta$  are selected as 8 and 3 respectively.

### **3.2 Identification of daily power consumption patterns in building operation**

As introduced in Section 2.1, the intrinsic periods in the time series of building total power consumption are estimated using the non-parametric spectral density estimation method. The top three dominant frequencies are 0.0103, 0.0417 and 0.1121, which correspond to periods of 97 (i.e.,  $1/0.0103$ ), 24 (i.e.,  $1/0.0417$ ) and 9 (i.e.,  $1/0.1121$ ) respectively. Since the BAS data are collected at an interval of 15-minute, these three periods are approximately 1-day, 6-hour and 2-hour respectively.

The dominant period in the sequence of building total power consumption is 1-day. Therefore, the whole BAS data are segmented into daily subsequences and then

1 transformed into SAX representations. Increasing the word size  $W$  and alphabet size  $A$   
2 will lead to a better SAX representation of the original time series. However, the  
3 reduction in computation load is less. Miller et al. recommended  $W$  and  $A$  as 4 and 3  
4 respectively to identify typical patterns in building power consumption data [32].  
5 Actually, the selection of  $W$  and  $A$  is influenced by the scale of the building,  
6 installation capacities (e.g., cooling, heating, total electricity power) and operation  
7 strategies. A large building with high installation capacities tends to require large  $W$   
8 and  $A$  to adequately describe the variation in the original time series data. In this  
9 study,  $W$  is chosen as 12, considering that 2-hour was identified as one of the  
10 dominant periods. Considering that the chiller plant usually accounts for a large  
11 proportion of the total power consumption and the maximum running chiller number  
12 is 5 in the BAS data to be analyzed,  $A$  is chosen as 5 to reflect there are five major  
13 levels of power consumption due to the on-off control of chillers. It should be noted  
14 that the standardization is only applied to the total building power consumption time  
15 series, but not the daily subsequences. The consideration here is to identify typical  
16 daily patterns considering both the shape and magnitude.

17 The SAX representations of daily subsequences are then partitioned into different  
18 groups using the EAC method.  $K_{lower}$  and  $K_{upper}$  are selected as 2 and 20 respectively.  
19 The iteration number  $E$  is set as 200. As a result, 8 clusters are identified. Clusters 5,  
20 6, 7 and 8 only consists of 6 daily subsequences out 365 subsequences. Those  
21 subsequences are actually subsequence-wise outliers, as their shape and magnitude  
22 are dramatically different from the others. They are excluded from further analysis.

Figure 3 presents the profiles of daily subsequences in Clusters 1 to 4. Further examination of each cluster shows that Clusters 1 to 4 can be best interpreted using the climate and day type. Cluster 1 includes weekends in cold season and Cluster 4 contains weekdays in hot season. Cluster 2 and Cluster 3 mainly include weekdays in cold season and weekends in hot season respectively. The clustering results are coincident with the results obtained in our previous study [4, 18] and domain knowledge. It indicates that the SAX transformation can very well preserve the important information in original time series data.

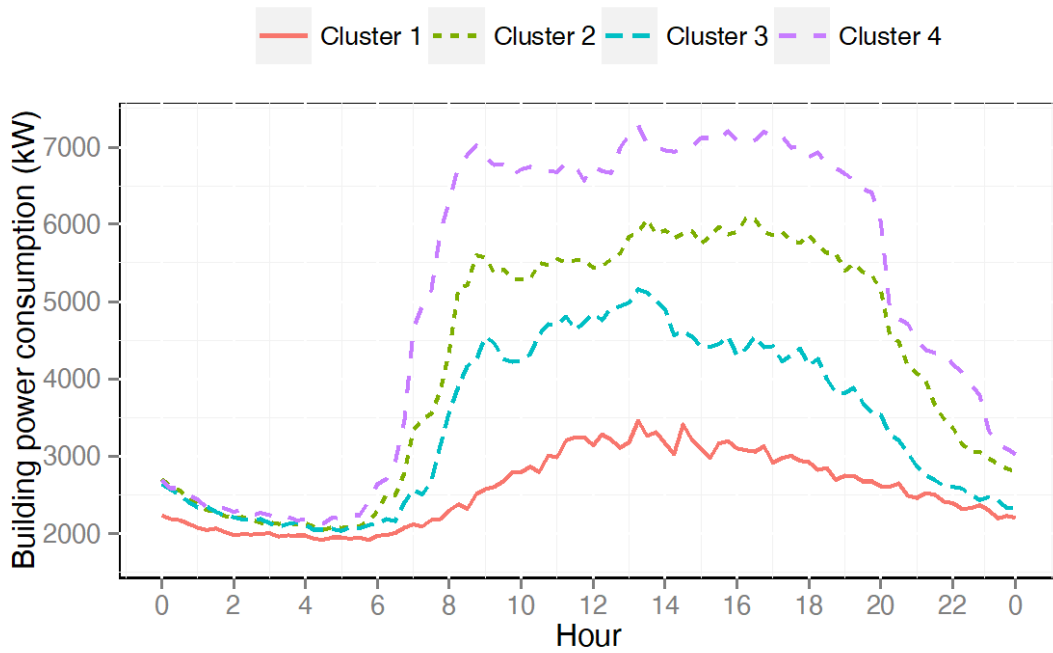


Figure 3: Four typical prototypes of daily building power consumption

### 3.3 Identify frequent operating patterns of subsystems

Univariate and multivariate motif discovery are applied to the 4 clusters separately to identify the frequent operating patterns. Considering that the daily operating

1 conditions (including outdoor weather conditions and indoor occupancy and  
2 equipment utilization conditions) varies largely, it is more meaningful to discover  
3 motifs in building operations with smaller lengths, compared with the above  
4 identification of power consumption pattern. In this study, the length of the univariate  
5 motifs to be discovered is set as 6-hour, as it is identified as the second dominant  
6 period in the building power consumption data. More specifically, subsequences are  
7 segmented using a 6-hour sliding window, which means the subsequences created are  
8 overlapping. Standardization is performed for each subsequence in each cluster. SAX  
9 representations are created using the setting of  $W=6$  and  $A=5$ . In such a case, each  
10 SAX symbol represents the hourly mean and has five possible levels. The iteration  
11 number for random projection is 100. During each iteration, 4 out of 6 SAX symbols  
12 are randomly selected for comparison, which means that subsequences belonging to  
13 the same motif can be different at one position at most [42].

14 Table 3 summarizes the number of univariate motifs discovered for each subsystem in  
15 Cluster 4 (i.e., weekdays in hot season). Figure 4 presents 4 motifs discovered in the  
16 time series of the aggregated chiller power consumption in Cluster 4. Each curve  
17 represents an occurrence of the corresponding motif. It is apparent that the time series  
18 subsequences belonging to the same motif are very similar in their shapes and  
19 magnitudes. An uptrend in chiller power consumption is observed in Figure 4a. It is  
20 shown that two chillers are sequentially switched on at the beginning of working  
21 hours (i.e., 6:00 a.m. to 9:00 a.m.) to cope with the upcoming morning peak of  
22 occupancy and equipment utilization. The chiller switch-off process shares a similar

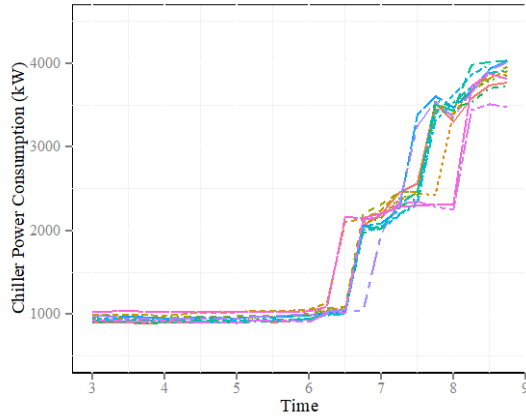
1 pattern and two chillers are sequentially switched of (Figure 4b). The other two  
2 motifs, as shown in Figures 4d and 4c, present relatively steady operating conditions.  
3 The chiller operation between 0:00 a.m. and 6:00 a.m. is steadily maintained at a low  
4 level due to the absence of occupancy. By contrast, the chiller power consumption is  
5 maintained at a much higher level between 9:00 a.m. to 3:00 p.m. A slight decrease  
6 can be observed from 1:00 p.m. to 2:00 p.m., which is in accordance with the lunch  
7 time for most companies in ICC.

8 Typical operating behaviors can be obtained by analyzing the univariate motifs  
9 identified. For instance, Figure 5 presents 2 frequent patterns for the AHU operation  
10 between 9:00 p.m. and 3:00 a.m. in Cluster 4. The main difference is that a sudden  
11 drop in AHU power consumption is observed at 12:00 a.m. in Figure 5a, while the  
12 AHU power consumption gradually decreases in Figure 5b. After carefully examined  
13 the original data, it is found that the AHU power consumption measured at three  
14 mechanical floors (i.e., 6/F, 42/F and 78/F) simultaneously drop at 12:00 a.m. in  
15 pattern 1. By contrast, the drops are observed at 10:00 p.m., 1:00 a.m. and 2:00 a.m.  
16 for the AHUs at 42/F, 78/F and 6/F slightly and gradually in pattern 2.

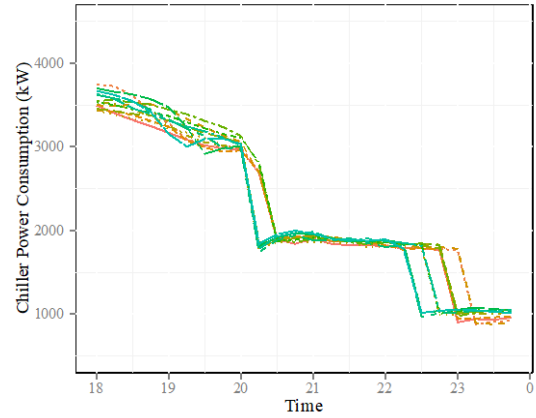
17 Table 3: A summary of univariate motifs discovered in Cluster 4

Subsystems	Chiller	CT	SCHWP	AHU	PAU	MV	VTS	NP	EP	PD
Motif No.	15	9	10	17	14	19	4	15	3	3

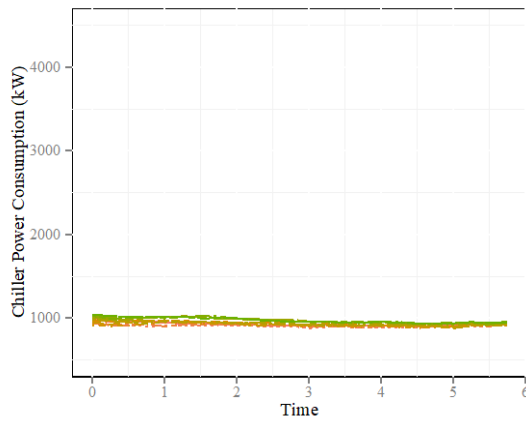
18



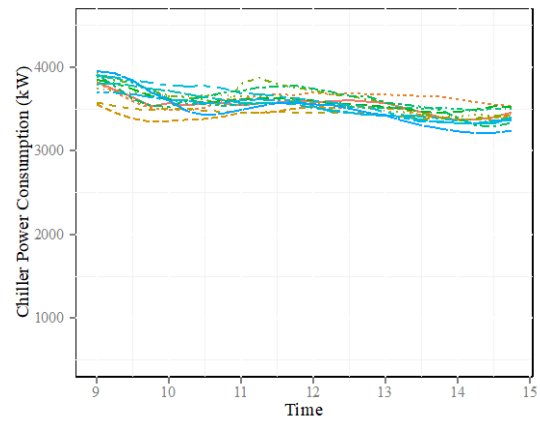
(a) Uptrend in chiller operation



(b) Downtrend in chiller operation

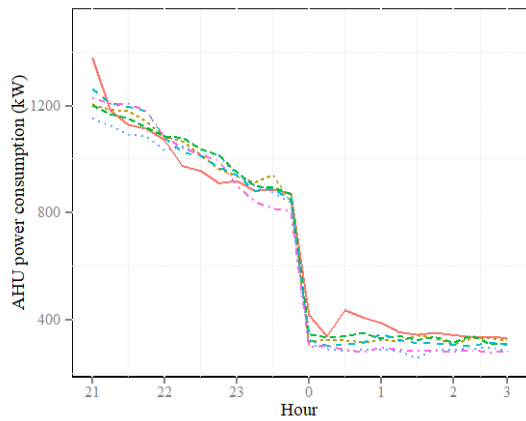


(c) Horizontal trend at low level

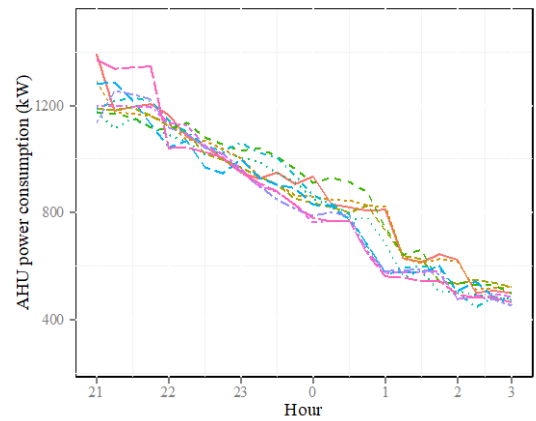


(d) Horizontal trend at high level

Figure 4: Examples of univariate motifs in chiller operation in Cluster 4



(a) Pattern 1



(b) Pattern 2

Figure 5: Typical AHU operating patterns between 9:00 p.m. to 3:00 a.m. in Cluster 4

As introduced in Section 2.3.1, univariate motifs discovered are used to identify

multivariate motifs. The algorithm can discover both synchronous and non-synchronous multivariate motifs. The parameter  $\alpha$  is set as 0.8. Figure 6 presents an example of simultaneous multivariate motif. It depicts the building dynamic operations for different subsystems during 3:00 a.m. to 9:00 a.m. The chiller power consumption starts to rise from 6:30 a.m. and two chillers are sequentially switched on. A rise in SCHWP power consumption can be observed accordingly to circulate the chilled water. The PAU power consumption stays steady at low-level until 8:00 a.m. This is because ICC adopts demand-controlled ventilation to control the PAU and the occupancy increases from 8:00 a.m. because people start to work. A rise in MV power consumption can be observed at 8:00 a.m., which is also due to the occupancy change. In addition, the MV power consumption also undergoes an increase at around 6:30 a.m., which relates to the activation of the precooling strategy. Similarly, uptrends in VTS and NLTG power consumptions can be observed in Figures 6e and 6f to cope with the increase in occupancy. These motifs show that the HVAC system in ICC is under reliable control and operations well meet the expectations. ICC was awarded as an Intelligent Building of 2011 by the Asian Institute of Intelligent Buildings, partly owing to the advanced BAS installed in ICC.

### **3.4 Identify temporal association rules between subsystem operations**

This section focuses on discovering the temporal associations between the operations of different subsystems. The operation of each subsystem at certain time instant is represented by two features, i.e., level and trend. The power consumption data of each

1 subsystem are categorized into three levels, i.e., *Low*, *Medium* and *High*. The trend is  
2 defined based on the changes between successive time step and categorized into 1 to  
3 7, indicating large decrease, moderate decrease, slight decrease, steady, slight  
4 increase, moderate increase and large increase. The categorization thresholds are  
5 determined using k-means clustering algorithm. The TRuleGrowth algorithm is  
6 applied with the minimum support and confidence being set as 0.2 and 0.8  
7 respectively. The maximum time span changes from 1 (i.e., 15-minute) to 12 (i.e.,  
8 3-hour). The post-mining method described in Section 2.4.2 is applied to find the  
9 exact time lag in temporal association rules.

10 Table 4 presents three example rules describing the inter-subsystem temporal  
11 associations in the multivariate motif shown in Figure 6. The first rule shows that  
12 when the AHU power consumption is *Low* and experiencing a slight increase at time  
13  $T$ , the chiller power consumption will be *Low* and stay steady at time  $T+1$ . The  
14 second rule shows that given the same antecedent, a slight increase in the chiller  
15 power consumption will be observed at  $T+2$ . These two rules demonstrate that the  
16 change in AHU and chiller operation is not synchronous and the time lag is around 15  
17 minutes. The last example rule describes the temporal association between the NLTG  
18 and the PAU power consumptions. It states that when the NLTG consumption is *Low*  
19 and experiencing a significant increase at time  $T$ , a significant increase in the PAU  
20 power consumption will be observed at  $T+9$ . The result's validity can be verified by  
21 manually inspecting Figure 6. For instance, the first significant increase in NLTG and  
22 PAU power consumptions take place at around 5:45 a.m. and 8:00 a.m. respectively

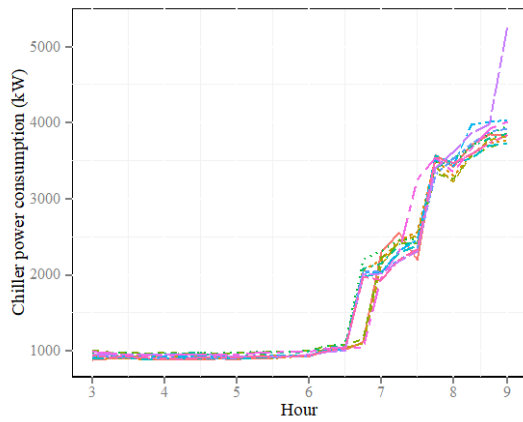


1 and therefore, the time lag for the third rule should be 9 unites of time (i.e., 135  
2 minutes).

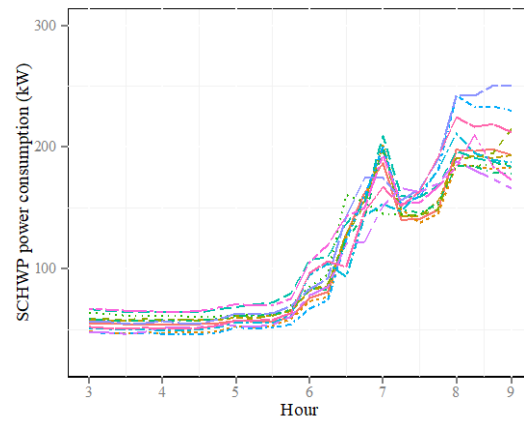
3 Table 4: Examples of temporal associations discovered

Rule	Antecedent	Consequent	Time lag (15-min)	Supp.	Conf.
1	AHU=Low, 5	Chiller=Low, 4	1	1.00	1.00
2	AHU=Low, 5	Chiller=Low, 5	2	0.89	0.89
3	NLTG=Low, 7	PAU=Low, 7	9	0.78	0.82

4

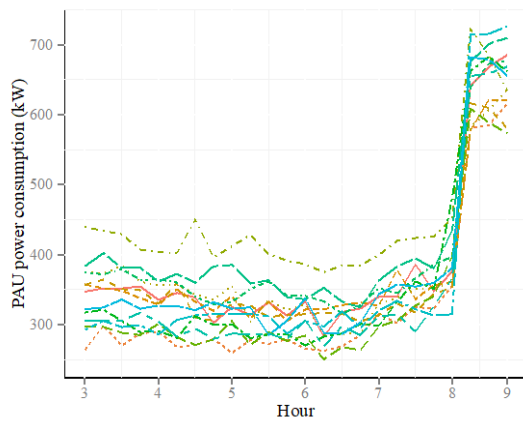


5

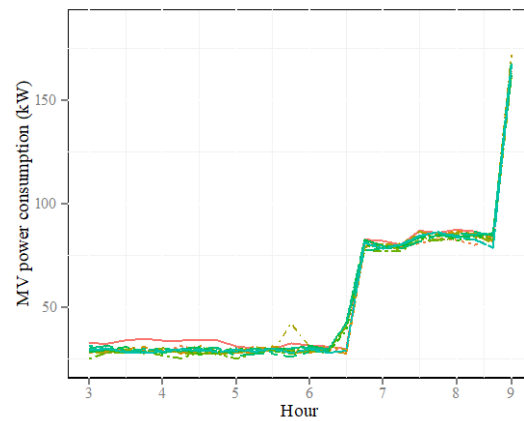


6 (a) Motif in chiller power consumption

(b) Motif in SCHWP power consumption

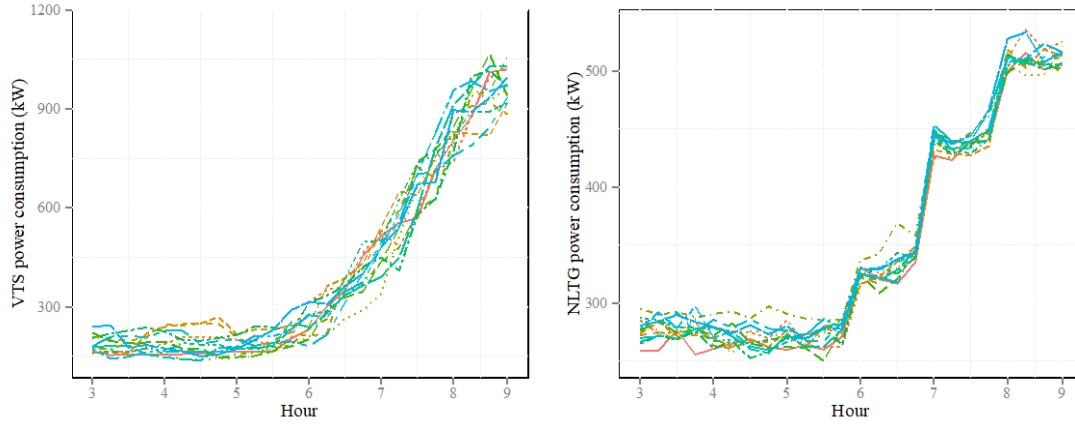


7



8 (c) Motif in PAU power consumption

(d) Motif in MV power consumption



(e) Motif in VTS power consumption (f) Motif in NLTG power consumption

Figure 6: An example of multivariate motif in Cluster 4

#### 4. Applications of Temporal Knowledge Discovered

A straightforward approach to applying the temporal knowledge discovered to building management is to build a database of motifs and temporal association rules as the benchmark of building operations. Then, the real-time BAS time series data are compared with the benchmarked operations to identify any possible anomalies. The post-mining methods developed in this study provide two more approaches to such applications. The following parts demonstrate these applications.

##### 4.1 Applications of associations between univariate motifs

The post-mining method introduced in Section 2.4.1 is applied to discover associations between univariate motifs. To illustrate, 103 univariate motifs which are discovered in Cluster 4 are used for analysis. The Apriori algorithm is applied with the minimum support and confidence set as 0.1 and 0.8 respectively. These thresholds are set in such a way to ensure the discovery of strong but not necessarily frequent

1 associations. 144 association rules are discovered. The association rules obtained can  
2 be applied to find anomalies in operation, such as less energy-efficient operations,  
3 faulty operations, as well as normal but rare operations.

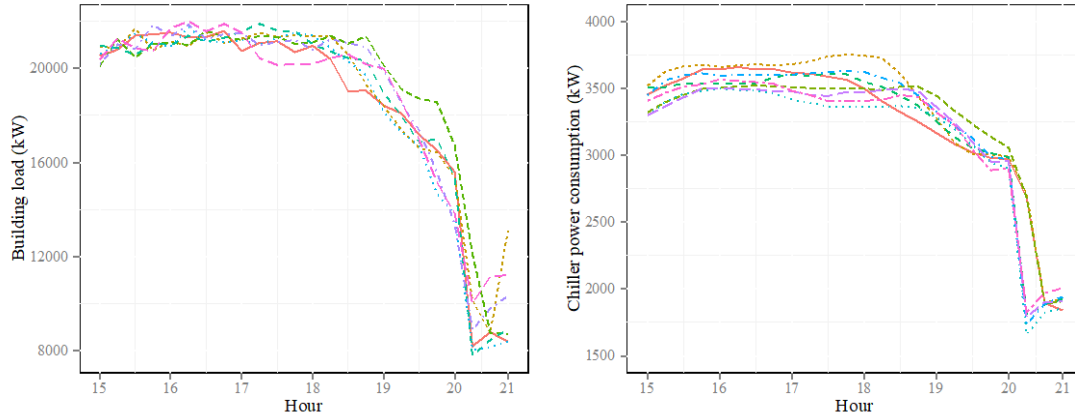
4 As shown in Figure 7, one rule is *Cooling Load = Motif 3*  $\rightarrow$  *Chiller = Motif 11*. It  
5 describes the association between Motif 3 in the building cooling load and Motif 11 in  
6 the chiller power consumption, which both take place between 3:00 p.m. to 9:00 p.m.

7 Atypical patterns are identified by finding the time series data which meet the  
8 antecedent but not the consequent. An example is presented in Figure 8. Motif 11 in  
9 the chiller power consumption is shown using blue boxplots and the atypical chiller

10 operation is shown using the red solid line. Given the same building load demand, the  
11 atypical operation results in much higher chiller power consumption during the period  
12 from 3:00 p.m. to 7:30 p.m. The mean chiller coefficient of performance (COP)

13 decreases from 5.82 to 5.12 (i.e. 12% drop in energy efficiency) when the atypical  
14 operation takes place. It is found out by examining original data that during chiller  
15 Motif 11, three chillers are running at a nearly full-load condition. By contrast, 4

16 chillers are switched-on during the atypical operation with a lower part-load ratio. In  
17 such a case, the identified atypical operation resulted in a less energy efficient  
18 operation.



(a) Antecedent motif

(b) Consequent motif

Figure 7: Association between building cooling load and chiller motifs in Cluster 4

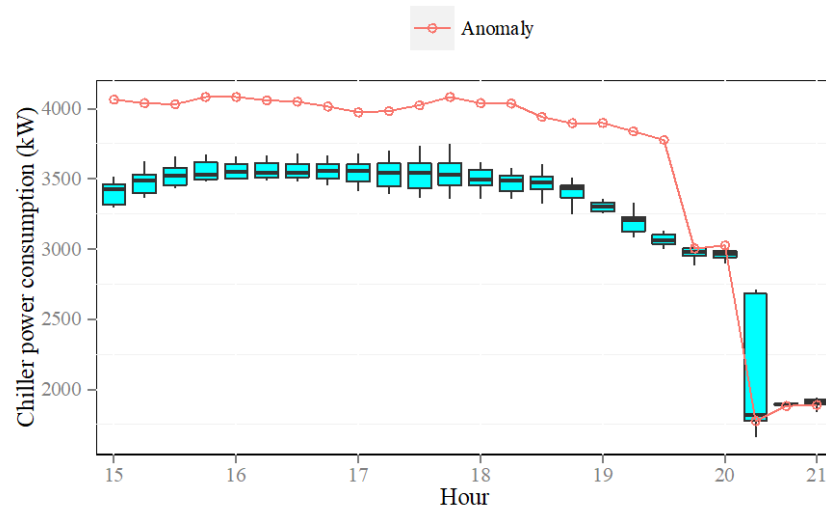
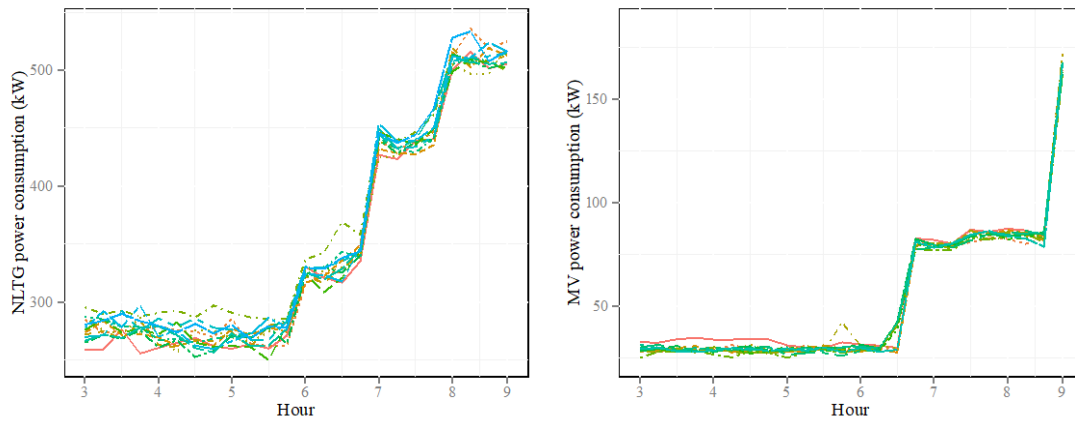


Figure 8: Comparison of chiller operations

Another example rule is  $NLTG = \text{Motif } 6 \rightarrow MV = \text{Motif } 9$ . It describes the association in the operating patterns of the normal power and lighting (NLTG) and mechanical ventilation (MV). These two motifs both take place between 3:00 a.m. to 9:00 a.m. and are shown in Figure 9. Figure 10 compares an atypical MV operation with the MV Motif 9. Starting from 4:30 a.m., the atypical operation has higher MV consumption than that in MV Motif 9. Further investigation shows that the difference

1 is caused by the MV at the third mechanical floor (i.e., 78/F). Normally, the MV  
2 consumption at the third mechanical floor is maintained at around 20kW between  
3 3:00 a.m. to 9:00 a.m. During the atypical operation, it experiences a sudden increase  
4 from 20kW to 45kW at 4:30 a.m. and is maintained at that level afterwards. One  
5 possible reason for this increase is due to the occupancy change in the corresponding  
6 office zone. However, the occupancy in office zone is unlikely to change at 4:30 a.m.  
7 In addition, the NLTG consumption is also subject to the influence of occupancy and  
8 no significant difference is observed during atypical operation. Such atypical  
9 operation may be due to the interference of manual control.



(a) Antecedent motif

(b) Consequent motif

Figure 9: Association between NLTG and MV motifs in Cluster 4

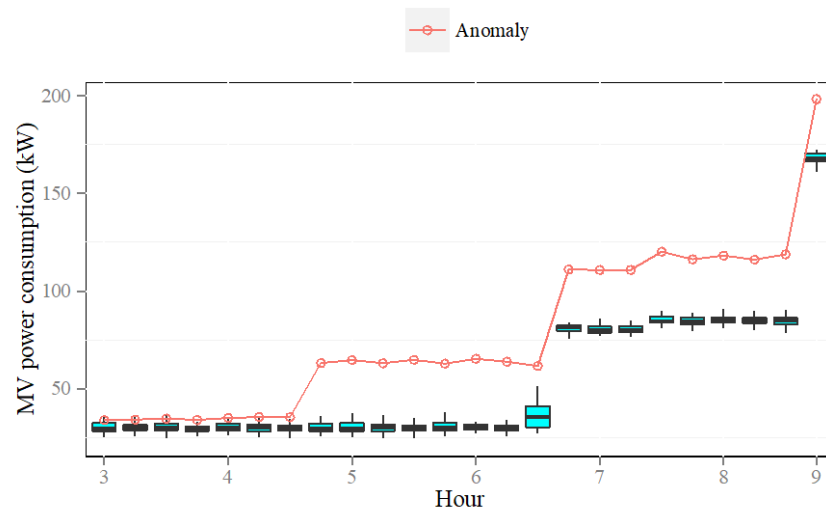
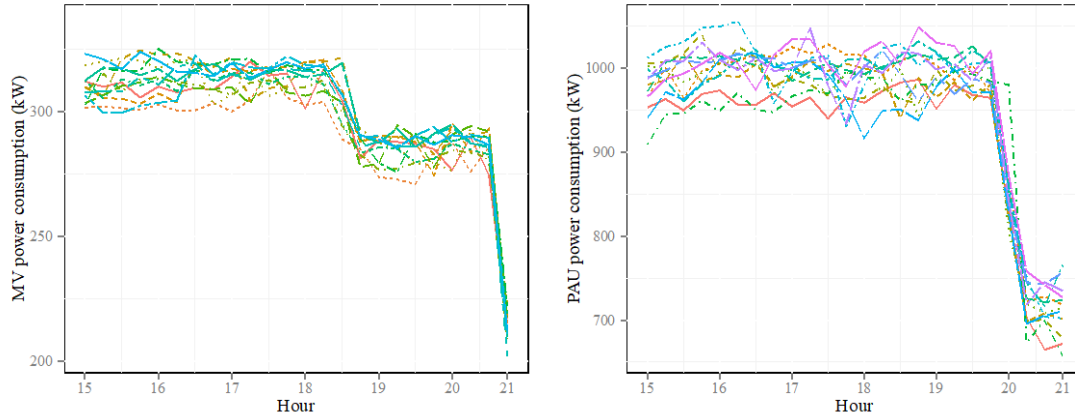


Figure 10: Comparison of MV operations

Another example rule describes the association between MV Motif 18 and PAU Motif 14. As shown in Figure 11, both motifs take place between 3:00 p.m. to 9:00 p.m. The two drops in MV consumption at around 6:30 p.m. and 8:45 p.m. are due to the decrease in MV consumption at the second and the first mechanical floors respectively. By contrast, one significant drop in the PAU consumption is observed at around 8:00 p.m., which is due to the huge decrease in office occupancy. An atypical operation is identified and its PAU consumption is compared with the PAU Motif 14 in Figure 12. Compared with PAU Motif 14, the PAU consumption in atypical operation is much smaller from 5:30 p.m. to 8:00 p.m. The reason behind is that the next day is a public holiday in Hong Kong and many offices have their employees released at around 5:00 p.m. Consequently, a power reduction in PAU consumption is observed. In such a case, the atypical operation identified is a normal but rare operation.



(a) Antecedent motif

(b) Consequent motif

Figure 11: Association between MV and PAU motifs in Cluster 4

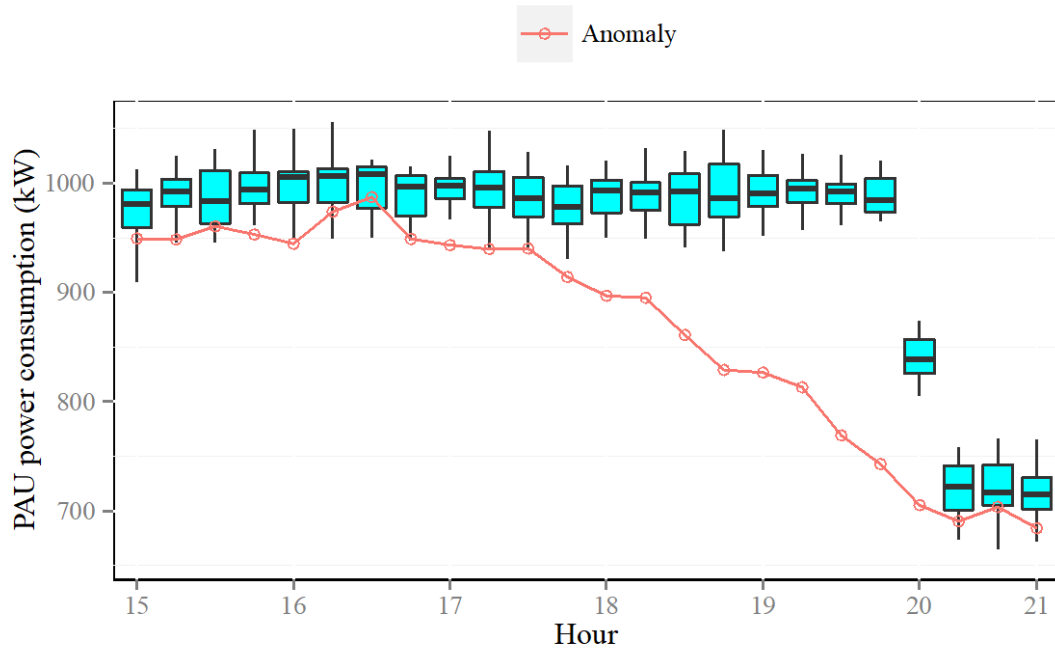


Figure 12: Comparison of PAU operations

## 4.2 Application of temporal association rules

### 4.2.1 Temporal anomaly detection

Temporal anomaly can be detected using the temporal association rules. Two approaches are possible. If the anomaly widely exists in the time series data, a

1 temporal association rule specifying such atypical association will be derived. In such  
2 a case, temporal anomaly can be detected by finding those observations which are in  
3 accordance with these temporal association rules. However, those anomaly data are  
4 seldom available. The second approach is more practically feasible. A knowledge  
5 database of normal temporal association rules can be constructed. Temporal anomaly  
6 can be detected by finding those observations which fail to meet the rules in the  
7 database.

8 An example is given here. Two rules with a time lag of 15-minute are derived to  
9 describe temporal associations in the chiller operation between 6:00 a.m. to

10 12:00p.m.:  $Chiller=High, 5 \xrightarrow{T=1} Chiller=High, 4$  and  $Chiller=High, 5 \xrightarrow{T=1}$

11  $Chiller=High, 3$ . These two rules specify that two possible operating modes are

12 possible at time  $T+1$  given the chiller power consumption at time  $T$  is *High* and has a

13 slightly increasing trend. Figure 13 presents the subsequences which fulfill these two

14 rules. The chiller power consumption at  $T+1$  will remain at *High* level, with either a

15 steady or a slightly decreasing trend. Temporal anomalies can be detected by finding

16 subsequences which fail to meet the rule consequent given the same antecedent.

17 Figure 14 presents an example of such anomalies. The anomaly is shown in red solid

18 line. It meets the rule antecedent at time  $T$ ; however, the operating mode at time  $T+1$

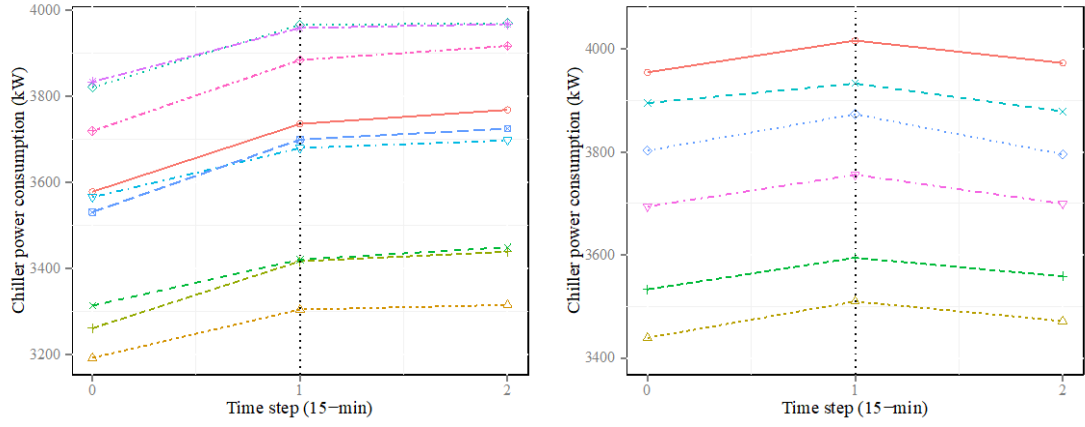
19 becomes *Medium* and has a significant decreasing trend. Further investigation shows

20 that at 8:30 a.m., Chiller 4 was switched off while two other chillers were switched on

21 as replacement. After consulting with the operation staff, it is found that Chiller 4 was

22 manually switched off due to its high operating current.





(a) Chiller=High, 5  $\xrightarrow{T=1}$  Chiller=High, 4      (b) Chiller=High, 5  $\xrightarrow{T=1}$  Chiller=High, 3

Figure 13: Examples of temporal associations in chiller operation

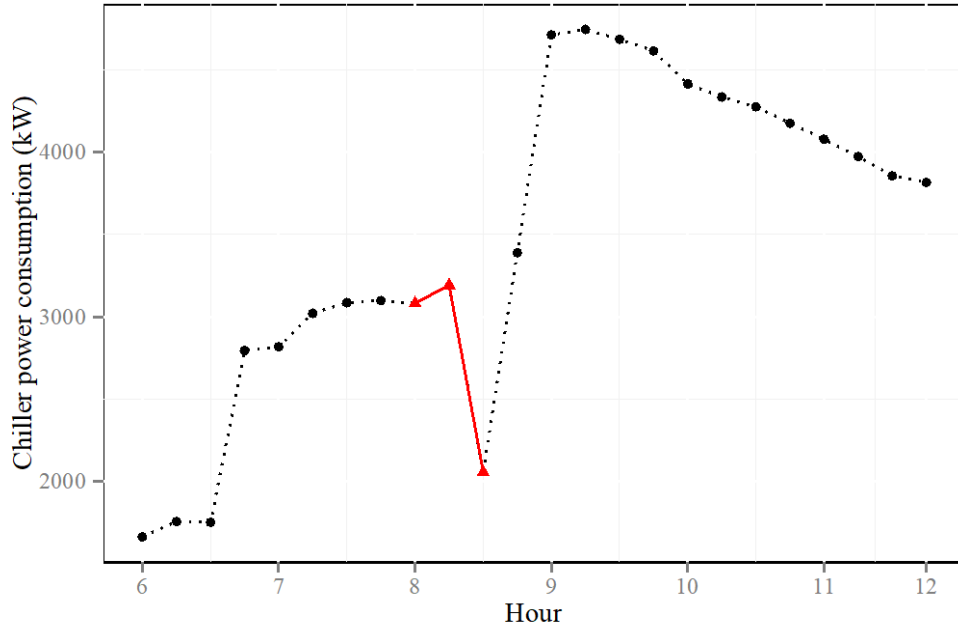


Figure 14: An example of temporal anomalies

#### 4.2.2 Characterization of building dynamics

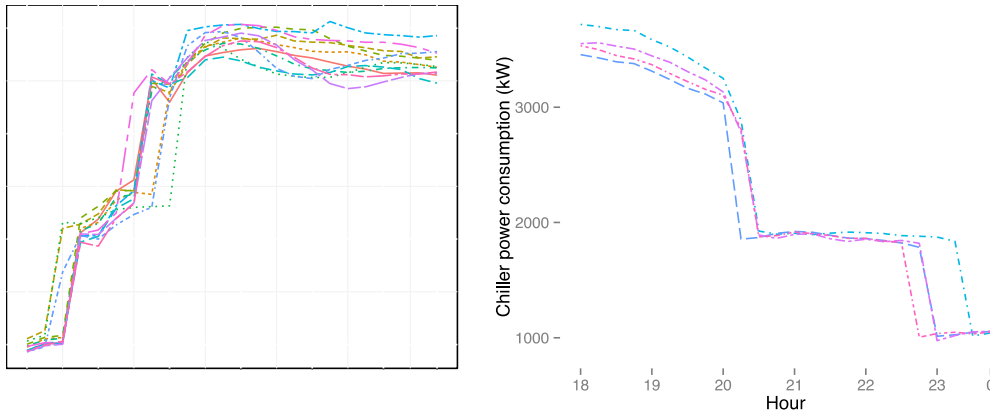
The extraction of time lag between the antecedent and the consequent of temporal association rules helps to characterize the building dynamics. Table 5 presents six example rules describing temporal associations in chiller operation. These rules are extracted from two chiller motifs, which are shown in Figure 15. The first three rules

1 are derived from the chiller Motif A, which takes place between 6:00 a.m. to 12:00  
2 p.m. Rule 1 indicates that if the chiller power consumption is *Low* and experiencing a  
3 significant increase at time  $T$ , the operating mode will be at *Medium* level and under  
4 slight increase at  $T+3$ . The second rule shows that once the chiller power  
5 consumption starts to increase significantly at *Low* level, it will reach its steady state  
6 at *High* level at  $T+10$ . Rule 3 shows that the time lag between two significant  
7 increases at *Low* and *Medium* levels is around 1 hour. The latter three rules are  
8 derived from the chiller Motif B, which occurs between 6:00 p.m. to 12:00 a.m. Rule  
9 4 states that if the chiller consumption is *High* and experiencing a significant decrease  
10 at time  $T$ , its steady state at *Medium* level will be reached at  $T+3$ , i.e., 45 minutes  
11 later. Similarly, Rule 5 describes that the time needed for the chiller power  
12 consumption to reach its steady state from the *Medium* to *Low* level is also 45  
13 minutes. The last rule quantifies that the time lag between the huge decrease at *High*  
14 and *Medium* levels is around 3 hours. The result is verified by checking Figure 15.  
15 The knowledge discovered in this subsection helps to quantify the building dynamics  
16 from two perspectives, i.e., the power consumption level and relative changes  
17 between successive time steps (i.e., trend). The temporal interactions and dynamics  
18 can be automatically extracted. Useful insights can be gained into how building  
19 subsystems react to a certain change in operation over time. The temporal associations  
20 discovered can be used to facilitate the optimal control and decision-makings in  
21 building operation, e.g., chiller sequence control and integration between individual  
22 buildings and large power grid systems.

1

Table 5: Temporal associations in chiller operations

Rule	Motif	Antecedent	Consequent	Time lag (15-minute per unit)	Support	Confidence
1	A	Chiller=Low, 7	Chiller=Medium, 5	3	0.92	0.97
2	A	Chiller=Low, 7	Chiller=High, 4	10	0.81	0.92
3	A	Chiller=Low, 7	Chiller=Medium, 7	4	0.83	0.83
4	B	Chiller=High, 1	Chiller=Medium, 4	3	0.36	0.87
5	B	Chiller=Medium, 1	Chiller=Low, 4	3	0.78	0.85
6	B	Chiller=High, 1	Chiller=Medium, 1	12	0.44	0.84



2

3 (a) Motif A: Between 6 a.m. to 12 p.m. (b) Motif B: Between 6 p.m. to 12 a.m.

4 Figure 15: Two examples of chiller operation motifs

5

6 **5. Conclusions and Discussions**

7 BAS data are in essence multivariate time series data. Currently, few studies have

8 addressed temporal knowledge discovery and applications in big BAS data. This

9 study proposes a generic temporal knowledge discovery methodology for mining big

10 BAS data. A diversity of time series data mining techniques and their practical

11 potentials in analyzing big BAS data for building operations and performance

12 management are explored in this study. Rather than addressing pre-defined specific

1 problems, the methodology developed mainly aims to discover unknown temporal  
2 knowledge by adopting unsupervised DM techniques to mine the big BAS data. The  
3 intention is to let the data tell the story and then, using domain knowledge to interpret,  
4 select and apply the knowledge discovered. The methodology proposed serves as a  
5 prototype of big data analysis tools which can be integrated with modern building  
6 automation systems to realize automatic knowledge discovery and applications.

7 This study specifically addresses two major challenges in mining big BAS data. One  
8 major challenge is the heavy computational load caused by the massive data amount.  
9 From a technological perspective, this challenge can be tackled by using  
10 high-performance computing machines or cloud-based computing. The adoption of  
11 suitable data transformation methods and more computationally efficient DM  
12 algorithms can provide an alternative solution. This study shows that the SAX method  
13 is capable of reducing the data numerosity while preserving the majority of the  
14 information contained in the BAS power consumption data. The univariate motif  
15 discovery algorithm adopted in this study is based on the concept of combinatorial  
16 search rather than exhaustive search and thereby the required computational costs can  
17 be largely reduced. Another challenge is the extraction of new features based on the  
18 original data for knowledge discovery, also known as feature engineering. Extraction  
19 of novel and unique features can greatly enhance the mining result quality. Besides  
20 the power consumption level, this study makes use of the changing trend to describe  
21 the mode of each subsystem at each time step. The temporal association rules  
22 discovered are more meaningful and straightforward for knowledge interpretation and

1 application, compared with those obtained using other features as inputs (e.g., the  
2 power consumption level alone).

3 Time series data mining can discover large amounts of knowledge with different  
4 types, such as clusters, univariate and multivariate motifs, and temporal association  
5 rules. It is challenging and time-consuming to interpret and apply the knowledge  
6 discovered. This study develops two methods for the efficient post-processing of  
7 knowledge discovered. The first method uses a co-occurrence matrix to map the  
8 relationship between univariate motifs. Reliable associations between univariate  
9 motifs are derived which provides a novel and convenient approach to utilizing  
10 univariate motifs. The second method utilizes a filtering method to improve the  
11 temporal association rules mining algorithms with the accurate estimation of time  
12 interval between the antecedent and the consequent. The time interval or lag provides  
13 valuable insights into building dynamics and HVAC performance characteristics. The  
14 methodology has been applied to analyze the BAS data retrieved from the tallest  
15 building in Hong Kong. The knowledge discovered has been successfully used to  
16 identify anomalies in building operations and characterize the building dynamics. The  
17 open-source software *R* and *SPMF* were used to perform the mining.

18 Enabling the building automation industry to benefit from advances in big data  
19 analysis is a non-trivial task. It requires building professionals to thoroughly  
20 understand the mechanisms of both DM algorithms and building operations. A lot  
21 more work needs to be done other than simply applying DM algorithms to analyze  
22 BAS data. The main purpose of this study is to bridge the knowledge gap between

building professionals and advanced data analytics. One limitation of this research is that it only considers power consumption data. The data transformation methods considering the physical variables in BAS data, such as the temperature, relative humidity, water flow rate and pressure, are more challenging. In addition, even though two methods have been developed to enhance the post-mining efficiency, the amount of knowledge to be examined by domain expertise is still large. For instance, 103 association rules are obtained from mining the univariate motifs in one data cluster. Future research will focus on developing transformation methods for various types of physical variables in BAS data and propose solutions to further enhance the post-mining efficiency.

## **Acknowledgements**

The authors gratefully acknowledge the support of this research by the Research Grant Council (RGC) of the Hong Kong SAR (152181/14E). Henrik Madsen is partly funded by the Danish CITIES project (DSF-1305-00027B) which is hence also acknowledged.

## **References**

- [1] L. Perez-Lombard, J. Ortiz, C. Pout, A review on buildings energy consumption information, *Energy and Buildings* 40 (2008) 394-398.
- [2] United Nations Environment Programme (UNEP), Common carbon metric for measuring energy use and reporting greenhouse gas emissions from building

1 operations, <http://www.unep.org/sbci/pdfs/UNEPSBCICarbonMetric.pdf>, accessed on  
2 April 9, 2015.

3 [3] P. Waide, J. Ure, N. Karagianni, G. Smith, B. Bordass, The scope for energy and  
4 CO<sub>2</sub> savings in the EU through the use of building automation technology, Final  
5 Report for the European Copper Institute, August 10, 2013.

6 [4] C. Fan, F. Xiao, C.C. Yan, A framework for knowledge discovery in massive  
7 building automation data and its applications in building diagnostics, *Automation in*  
8 *Construction* 50 (2014) 81-90.

9 [5] F. Dalene, Technology and information management for low-carbon building,  
10 *Journal of Renewable and Sustainable Energy* 4-041402 (2012),  
11 doi:10.1063/1.3694120.

12 [6] B. Dong, C. Cao, S.E. Lee, Applying support vector machines to predict building  
13 energy consumption in tropical region, *Energy and Buildings* 37 (2005) 545-553.

14 [7] A.S. Ahmad, M.Y. Hassan, M.P. Abdullah, H.A. Rahman, F. Hussin, H. Abdullah,  
15 R. Saidur, A review on applications of ANN and SVM for building electrical energy  
16 consumption forecasting, *Renewable and Sustainable Energy Reviews* 33 (2014)  
17 102-109.

18 [8] C. Fan, F. Xiao, S.W. Wang, Development of prediction models for next-day  
19 building energy consumption and peak power demand using data mining techniques,  
20 *Applied Energy* 127 (2014) 1-10.

21 [9] A. Kusiak, G.L. Xu, Z.J. Zhang, Minimization of energy consumption in HVAC  
22 systems with data-driven models and an interior-point method, *Energy Conversion*

1 and Management 85 (2014) 146-153.

2 [10] S.K. Kwok, K.K. Yuen, W.M. Lee, An intelligent approach to assessing the  
3 effect of building occupancy on building cooling load prediction, Building and  
4 Environment 46 (2011) 1681-1690.

5 [11] S.M. Wu, D. Clements-Croome, Understanding the indoor environment through  
6 mining sensory data – A case study, Energy and Buildings 39 (2007) 1183-1191.

7 [12] G. Kim, L. Schaefer, T.S. Lim, J.T. Kim, Thermal comfort prediction of an  
8 underfloor air distribution system in a large indoor environment, Energy and  
9 Buildings 64 (2013) 323-331.

10 [13] Z. Yu, F. Haghighat, C.M. Fung, H. Yoshino, A decision tree method for  
11 building energy demand modeling, Energy and Buildings 42 (2010) 1637-1646.

12 [14] J.S. Chou, Y.C. Hsu, L.T. Lin, Smart meter monitoring and data mining  
13 techniques for predicting refrigeration system performance, Expert Systems with  
14 Applications 41 (2014) 2144-2156.

15 [15] E.U. Kucuksille, R. Selbas, A. Sencan, Prediction of thermodynamic properties  
16 of refrigerants using data mining, Energy Conversion and Management 52 (2011)  
17 836-848.

18 [16] D.F.M. Cabrera, H. Zareipour, Data association mining for identifying lighting  
19 energy waste patterns in educational institutes, Energy and Buildings 62 (2013)  
20 210-216.

21 [17] A. Capozzoli, F. Lauro, I. Khan, Fault detection analysis using data mining  
22 techniques for a cluster of smart office buildings, Expert Systems with Applications



1 42 (2015) 4324-4338.

2 [18] F. Xiao, C. Fan, Data mining in building automation system for improving  
3 building operational performance, *Energy and Buildings* 75 (2014) 109-118.

4 [19] X. Xue, S.W. Wang, Y.J. Sun, F. Xiao, An interactive building power demand  
5 management strategy for facilitating smart grid optimization, *Applied Energy* 116  
6 (2014) 297-310.

7 [20] A. Azadeh, M. Saberi, S.F. Ghaderi, A. Gitiforouz, V. Ebrahimipour, Improved  
8 estimation of electricity demand function by integration of fuzzy system and data  
9 mining approach, *Energy Conversion and Management* 49 (2008) 2165-2177.

10 [21] I. Fernandez, C. Borges, Y. Penya, Efficient building load forecasting, In  
11 *Proceedings of the 16<sup>th</sup> IEEE International Conference of Emerging Technologies and*  
12 *Factory Automation*, 5-9 September, 2011, Toulouse, France.

13 [22] Y. Yao, Z.W. Lian, S.Q. Liu, Z.J. Hou, Hourly cooling load prediction by a  
14 combined forecasting model based on analytic hierarchy process, *International*  
15 *Journal of Thermal Sciences* 43 (2004) 1107-1118.

16 [23] M. Kawashima, C.E. Dorgan, J.W. Mitchell, Hourly thermal load prediction for  
17 the next 24h by ARIMA, EWMA, LR, and an artificial neural network, *ASHRAE*  
18 *Transactions* 101 (1995) 186-200.

19 [24] J.C.M. Yiu, S.W. Wang, A multiple ARMAX modeling scheme for forecasting  
20 of air conditioning system performance, *Energy Conversion and Management* 48  
21 (2007) 2276-2285.

22 [25] F. Zamora-Martinez, P. Romeu, P. Botella-Rocamora, J. Pardo, Towards energy

1 efficiency: forecasting indoor temperature via multivariate analysis, *Energies* 6 (2013)

2 4639-4659.

3 [26] L. Renners, R. Bruns, J. Dunkel, Situation-aware energy control by combining

4 simple sensors and complex event processing, *Workshop on AI Problems and*

5 *Approaches for Intelligent Environment*, August 2012, Montpellier, France.

6 [27] Y.C. Wen, G.Y. Lin, T. Sung, M. Liang, G. Tsai, M.W. Feng, A complex event

7 processing architecture for energy and operation management, *The 5<sup>th</sup> ACM*

8 *International Conference on Distributed Event-Based Systems*, July 11-15, 2011, New

9 York, USA.

10 [28] J. O'Donnell, E. Corry, S. Hasan, M. Keane, E. Curry, Building performance

11 optimization using cross-domain scenario modeling, linked data, and complex event

12 processing, *Building and Environment* 62 (2013) 102-111.

13 [29] T.C. Fu, A review on time series data mining, *Engineering Applications of*

14 *Artificial Intelligence* 17 (2011) 164-181.

15 [30] H. Madsen, *Time series analysis*, 1<sup>st</sup> edition, Chapman & Hall/CRC Texts in

16 *Statistical Science*, 2007.

17 [31] D. Patnaik, M. Marwah, R.K. Sharma, N. Ramakrishnan, Temporal data mining

18 approaches for sustainable chiller management in data centers. *ACM Transactions on*

19 *Intelligent Systems and Technology* 2 (2011) 1-29.

20 [32] C. Miller, Z. Nagy, A. Schlueter. Automated daily pattern filtering of measured

21 building performance data, *Automation in Construction* 49 (2015) 1-17.

22 [33] A. Gelman, J. Hill, *Data analysis using regression and multi-level/hierarchical*

1 models, in: Analytical Methods for Social Research, Cambridge University Press, 1<sup>st</sup>  
2 edition, 2006.

3 [34] M. Gupta, J. Gao, C.C. Aggarwal, J.W. Han, Outlier detection for temporal data:  
4 A survey, IEEE Transactions on Knowledge and Data Engineering 15 (2014) 1-20.

5 [35] R.K. Pearson, Outliers in process modeling and identification, IEEE Transactions  
6 on Control Systems Technology 10 (2002) 55-63.

7 [36] J. Lin, E. Keogh, L. Wei, S. Lonardi, Experiencing SAX: A novel symbolic  
8 representation of time series, Data Mining and Knowledge Discovery 15 (2007)  
9 107-144.

10 [37] J. Kwac, J. Flora, R. Rajagopal, Household energy consumption segmentation  
11 using hourly data. IEEE Transactions on Smart Grid 5 (2014) 420-430.

12 [38] R. Gulbinas, A. Khosrowpour, J. Taylor, Segmentation and classification of  
13 commercial building occupants by energy-use efficiency and predictability, IEEE  
14 Transactions on Smart Grid 6 (2015) 1414-1424.

15 [39] C.S. Daw, C.E.A. Finney, E.R. Tracy, A review of symbolic analysis of  
16 experimental data, Review of Scientific Instruments 74 (2003) 915-930.

17 [40] A.L.N. Fred, A.K. Jain, Combining multiple clustering using evidence  
18 accumulation, IEEE Transactions on Pattern Analysis and machine intelligence 27  
19 (2005) 835-850.

20 [41] S. Vega-Pons, J. Ruiz-Schulcoper, A survey of clustering ensemble algorithms,  
21 International Journal of Pattern Recognition and Artificial Intelligence 25 (2011)  
22 337-372.

- 1 [42] B. Chiu, E. Keogh, S. Lonardi, Probabilistic discovery of time series motifs,  
2 ACM SIGKDD, Washington, DC, USA, 2003, pp. 493-498.
- 3 [43] Y. Tanaka, K. Iwamoto, K. Uehara, Discovery of time-series motif from  
4 multi-dimensional data based on MDL principle, Machine Learning 58 (2005)  
5 269-300.
- 6 [44] D. Minnen, C.L. Isbell, I. Essa, T. Starner, Discovering multivariate motifs using  
7 subsequence density estimation and greedy mixture learning, The 22<sup>nd</sup> National  
8 Conference on Artificial Intelligence, Volumn 1, 2007, pp. 615-620.
- 9 [45] A. Vahdatpour, N. Amini, M. Sarrafzadeh, Towards unsupervised activity  
10 discovery using multi-dimensional motif detection in time series, IJCAI 2009 21<sup>st</sup>  
11 International Joint Conference on Artificial Intelligence.
- 12 [46] M.J. Zaki, SPADE: An efficient algorithm for mining frequent sequences,  
13 Machine Learning 42 (2001) 30-60.
- 14 [47] P. Fournier-Viger, U. Faghihi, R. Nkambou, E.M. Nguifo, CMRules: Mining  
15 sequential rules common to several sequences, Knowledge-Based Systems 25 (2012)  
16 63-76.
- 17 [48] P. Fournier-Viger, C.W. Wu, V.S. Tseng, R. Nkambou, Mining sequential rules  
18 common to several sequences with the window size constraint, Advances in  
19 Artificial-25<sup>th</sup> Canadian Conference on Artificial Intelligence, 2012, Toronto, Canada,  
20 pp. 299-304.