

1     **A variation focused cluster analysis strategy to identify typical daily heating**  
2                     **load profiles of higher education buildings**

3                     Zhenjun Ma<sup>a,\*</sup>, Rui Yan<sup>a,\*</sup>, Natasa Nord<sup>b</sup>

4                     <sup>a</sup>Sustainable Buildings Research Centre, University of Wollongong, NSW, 2522, Australia

5                     <sup>b</sup>Department of Energy and Process Engineering, Norwegian University of Science and  
6                                     Technology, Norway

7                     \*Email: [zhenjun@uow.edu.au](mailto:zhenjun@uow.edu.au); [ry721@uowmail.edu.au](mailto:ry721@uowmail.edu.au)

8     **Abstract:** This paper presents a variation focused cluster analysis strategy to identify typical  
9     daily heating energy usage profiles of higher education buildings. Different from the existing  
10    cluster analysis studies which were primarily developed using Euclidean distance as the  
11    dissimilarity measure and tended to group the daily load profiles with similar magnitudes,  
12    Partitioning Around Medoids (PAM) clustering algorithm with Pearson Correlation Coefficient-  
13    based dissimilarity measure was used in this study to group the daily load profiles on the basis of  
14    the variation similarity. A comparison of the proposed strategy with a k-means cluster analysis  
15    with Euclidean distance as the dissimilarity measure was also performed. The performance of the  
16    proposed strategy was tested and evaluated using the three-year hourly heating energy usage data  
17    collected from 19 higher education buildings in Norway. The results demonstrated the  
18    effectiveness of the proposed strategy in identifying the typical daily energy usage profiles. The  
19    identified typical heating load profiles provided the information such as the peaks and troughs of  
20    the daily heating demand, daily high heating demand period and daily load variation. The  
21    identified profiles also helped to categorize multiple buildings into different groups in terms of  
22    the similar energy usage behaviors to support further energy efficiency initiatives.

23 **Keywords:** Cluster analysis; Load profile; Pearson Correlation Coefficient; Higher education  
24 buildings

25 **Nomenclature**

26	$C$	set of the identified clusters
27	$cov$	covariance
28	$d$	distance
29	$D$	Dunn index
30	$k$	number of clusters
31	$n$	number of observations
32	$N_d$	number of days belongs to a typical daily load profile
33	$N_{d,max}$	maximum number of days belongs to a typical daily load profile
34	$o$	data point identified as a medoid
35	$p$	tail area probability
36	PCC	Pearson Correlation Coefficient
37	$q$	data point
38	$R$	studentized deviate
39	RP	relative proportion
40	$t$	t-distribution
41	$X, Y$	vectors
42	$x,y$	values of individual dimension
43	<b><i>Greek letters</i></b>	
44	$\alpha$	significance level
45	$\lambda$	critical value
46	$\sigma$	standard deviation

47  $\phi$  identified clusters

48 ***Subscripts***

49 ED Euclidean distance

50 PCC Pearson correlation coefficient

51

52 **1. Introduction**

53 Building energy efficiency is essential for reducing global energy usage and promoting  
54 environmental sustainability, as the building sector contributes to a large proportion of the total  
55 energy usage worldwide [1, 2]. With the development of automatic meter reading systems,  
56 massive high-resolution energy usage data from buildings can now be easily collected with a  
57 reasonably low cost [3]. This massive amount of data provides a great opportunity to assist in  
58 better understanding building energy usage characteristics and operational performance, and in  
59 extracting the useful and hidden information to support the areas including but not limited to  
60 building energy performance assessment and benchmarking, building load estimation and  
61 demand side management, occupant behavior prediction, and fault detection and diagnosis of  
62 heating, ventilation and air-conditioning systems.

63 Identification of typical building load profiles based on the collected massive energy usage  
64 data has been proved to be an effective way to understand building energy usage characteristics  
65 and help to develop cost effective load shifting and peak demand control strategies [4, 5]. Cluster  
66 analysis, as a data mining technique to discover the natural grouping(s) of a set of patterns, points,  
67 or objects [6], has been used in a number of studies to identify typical building load profiles [4, 5,  
68 7, 8]. Jota et al. [4], for instance, used an agglomerative hierarchical clustering algorithm with  
69 Euclidean distance (ED) to identify the typical building load profiles, which were further used to  
70 predict the accumulated energy usage at the end of the day and the daily peak demand. Typical

71 heating load profiles of Danish single-family detached homes were studied by do Carmo and  
72 Christensen [5] using the k-means algorithm. Three types of typical load profiles, i.e. high  
73 demand, medium demand and low demand, were identified for the buildings operated during  
74 weekdays and weekends, respectively. A binary regression analysis was also performed to  
75 identify the explanatory factors governing the different heating load profiles. The implementation  
76 and evaluation of a cluster analysis approach for smart meter data were reported by Flath et al.  
77 [7], in which the k-means algorithm was used to identify typical building daily and weekly load  
78 profiles of a business intelligence environment. Symbolic Aggregate approxXimation (SAX)  
79 method was used by Miller et al. [8] to transform building energy usage data into alphabets while  
80 the k-means algorithm was used to identify the typical daily load profiles. Fuzzy c-means (FCM)  
81 was adopted by Fernandes et al. [9] to identify the typical gas consumption profiles of residential  
82 buildings. It was found that the gas consumption peaks were related to the upper-middle social  
83 class with a high income and the highest daytime off-peak gas usage was related to the ageing  
84 population. Panapakidis et al. [10] utilized several clustering algorithms, including k-means, k-  
85 means++, minimum variance criteria, FCM and self-organizing map (SOM), to identify typical  
86 building electricity usage profiles. It was concluded that SOM and k-means++ in the frequency  
87 domain outperformed the other clustering techniques in terms of the clustering error.

88 Cluster analysis distinguishes data vectors based on a certain type of dissimilarity measures  
89 [11]. Although different cluster analysis algorithms have been proposed for different scenarios, to  
90 the best knowledge of the authors, the existing studies on the identification of building typical  
91 load profiles using cluster analysis were primarily developed using ED as the dissimilarity  
92 measure. Cluster analysis using ED-based dissimilarity measure tends to identify the daily load  
93 profiles that are similar in terms of the intensity rather than the variation. In other words, the  
94 typical daily load profile identified using cluster analysis with ED as the dissimilarity measure is

95 more related to the load magnitude. For example, do Carmo and Christensen [5] labeled the  
96 identified load profiles as high demand, medium demand, and low demand. ED-based  
97 dissimilarity measure is also difficult to identify building daily load profiles with similar  
98 variations but with different magnitudes, which will be elaborated in Section 2.2.

99 Higher education buildings have an important role in the minimization of greenhouse gas  
100 emissions from the built environment and in assisting the mitigation and adaptation of our society  
101 to climate change [12]. This paper presents a strategy using Partitioning Around Medoids (PAM)  
102 clustering algorithm to identify typical daily heating energy usage profiles of a group of higher  
103 education buildings. The novelty of this paper is to use Pearson Correlation Coefficient (PCC) as  
104 the dissimilarity measure to cluster daily heating energy usage profiles, in which the typical  
105 energy usage profiles are identified based on the load variation instead of the load magnitude,  
106 which is different from the majority of the previous studies used cluster analysis with ED as the  
107 dissimilarity measure. Based on the identified typical load profiles, a hierarchical clustering was  
108 used to group the buildings with similar heating energy usage characteristics. A comparison of  
109 the proposed strategy with an ED-based k-means cluster analysis strategy was also performed.  
110 The performance of the proposed strategy was evaluated using three-year hourly district heating  
111 energy usage data collected from 19 higher education buildings in Norway.

## 112 **2. Development of the variation focused cluster analysis strategy**

### 113 2.1 Outline of the variation focused cluster analysis strategy

114 The outline of the proposed variation focused cluster analysis strategy is illustrated in Fig. 1,  
115 which was developed following the standard Knowledge Discovery from Database (KDD)  
116 process [13]. It mainly consisted of four steps, including data collection, data pre-processing, data  
117 mining, and results evaluation and interpretation.

118 The collection of hourly energy usage data of individual buildings was the first step and the  
119 necessary data can be generally collected from building management systems. There were four  
120 tasks in the data pre-processing step, including outlier removal, data standardization, data  
121 segmentation and the removal of the weekend data and the data segments with small variations.  
122 In this study, the generalized Extreme Studentized Deviate (ESD) test method was used to  
123 identify and remove the outliers in the collected raw data. As the magnitude of the energy usage  
124 varied from building to building, to avoid the influence of identifying typical daily energy usage  
125 profiles, the processed data of each building was standardized to zero mean and one standard  
126 deviation. Data segmentation was then performed to transform the data into 24 hours segments in  
127 order to form daily load profiles. As the primary focus of this strategy was to identify the typical  
128 daily energy usage profiles during the building occupied periods with distinctive variation  
129 patterns, the segments during the weekends and the segments with small variations were  
130 discarded. The segments with small variations refer to the segments with a small difference  
131 between the daily maximum and minimum energy usages. In this study, a threshold of 5.0% was  
132 used, which means that 5.0% of the segments with the least difference among all daily segments  
133 were discarded.

134 In the data mining step (see Fig. 1), Pearson Correlation Coefficient (PCC) was first  
135 calculated to measure the dissimilarities among different daily load profiles. The Partitioning  
136 Around Medoids (PAM) clustering algorithm was then applied to cluster the daily load profiles  
137 with similar variations based on the PCC-based dissimilarity measure calculated. A boxplot was  
138 used to remove the daily load profiles with the large aggregated dissimilarities (i.e. the sum of the  
139 dissimilarities to all other daily load profiles in the same cluster) in each cluster, in order to  
140 reduce the influence of the extreme daily load profiles on the identification of typical daily load  
141 profiles. The daily load profiles with the aggregated dissimilarity measure beyond  $Q3+1.5IQR$ ,

142 where Q3 is the third quartile and IQR is an inter-quartile range between Q1 and Q3, were  
 143 discarded. The typical daily load profiles were then determined by averaging the remaining daily  
 144 load profiles in each cluster. Lastly, a hierarchical clustering was used to group the buildings with  
 145 similar load characteristics.

146 In the last step, the identified typical daily load profiles and building groups were visualized,  
 147 evaluated and interpreted.

## 148 2.2 Outlier removal with the generalized Extreme Studentized Deviate (ESD) test method

149 Generalized ESD test method has been applied for identifying and removing outliers in  
 150 building energy usage data in a number of studies [14-16]. This method detects outliers through  
 151 comparing the studentized deviate  $R$  of  $n$  extreme observations to a critical value  $\lambda$ . The extreme  
 152 observations are the observations with the first  $n$  largest differences compared to the mean value  
 153  $\bar{x}$ . The  $R_i$  of the  $i^{th}$  extreme observation  $x_{e,i}$  is determined using Eq. (1) and the corresponding  $\lambda_i$   
 154 is defined in Eq. (2) [14]. The generalized ESD test method starts with the most extreme  
 155 observation and compares its  $R_i$  to the corresponding  $\lambda_i$ . If  $R_i$  is greater than  $\lambda_i$ , the extreme  
 156 observation is then identified as an outlier and removed from the dataset. The same process is  
 157 applied to the next extreme observation until all the  $n$  extreme observations are examined. More  
 158 details of the generalized ESD test method can be found in [14]. If an outlier is identified and  
 159 removed, its position will be filled through the linear interpolation.

$$160 \quad R_i = \frac{|x_{e,i} - \bar{x}|}{\sigma} \quad (1)$$

$$161 \quad \lambda_i = \frac{(n-i)t_{n-i-1,p}}{\sqrt{(n-i+1)(n-i-1+t_{n-1-1,p}^2)}} \quad (2)$$

162 where  $\sigma$  is the standard deviation,  $t_{n-i-1,p}$  is the t-distribution with  $n-i-1$  degrees of freedom and  $p$  is  
 163 the tail area probability and is defined in Eq. (3) [14].

164 
$$p = \frac{\alpha}{2(n-i+1)} \quad (3)$$

165 where  $\alpha$  is the significance level.

### 166 2.3 Pearson Correlation Coefficient (PCC)-based dissimilarity measure

167 Cluster analysis groups the data by minimizing the inter-cluster dissimilarity while  
 168 maximizing the intra-clusters based on a certain type of the dissimilarity measures [17]. In the  
 169 proposed strategy, the distance between the two daily load profiles ( $d_{PCC}$ ) determined by Eq. (4)  
 170 was used to measure the dissimilarity between the two daily load profiles, in which the PCC is  
 171 defined in Eq. (5).

172 
$$d_{PCC}(X, Y) = 1 - PCC \quad (4)$$

173 
$$PCC(X, Y) = \frac{cov(X, Y)}{\sigma_X \sigma_Y} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}} \quad (5)$$

174 where  $d$  means the distance,  $cov$  stands for the covariance,  $X$  and  $Y$  represent the vectors, and  $x$   
 175 and  $y$  stands for the values of the individual dimension.

176 A comparison between the use of the PCC-based and ED-based dissimilarity measures is  
 177 illustrated in Fig. 2, where ED was calculated using Eq. (6). The data used in Fig. 2 was given  
 178 only for illustration purpose.

179 
$$d_{ED}(X, Y) = \sqrt{\sum_{i=1}^n (x_i - y_i)^2} \quad (6)$$

180 It can be seen that the ED of Profiles 1 and 2 ( $d_{12}$ ) and ED of Profiles 1 and 3 ( $d_{13}$ ) were  
 181 38.91 kWh and 8.178 kWh, respectively. Compared to Profile 2, Profile 3 was closer to Profile 1  
 182 in terms of the ED dissimilarity measure. However, the variation of Profile 2 was more similar to  
 183 that of Profile 1, as shown in Fig. 2(a). The PCC of Profiles 1 and 2 ( $PCC_{12}$ ) and PCC of Profiles



184 1 and 3 ( $PCC_{13}$ ) were 0.978 and 0.173, respectively. A higher PCC indicated a higher similarity  
 185 between the two profiles in terms of the daily load variation. Therefore, PCC-based dissimilarity  
 186 measure can better identify the daily load profiles with similar variations.

#### 187 2.4 Partitioning Around Medoids clustering algorithm

188 Partitioning Around Medoids (PAM) clustering algorithm [18] was used to cluster daily load  
 189 profiles using the PCC-based dissimilarity measure. In PAM, a medoid is a data point in a  
 190 particular cluster which has a minimized aggregated distance to all other data points in that  
 191 cluster. The objective of PAM clustering algorithm is to find a subset  $\{o_1, o_2, \dots, o_k\} \in \{q_1, q_2, \dots, q_n\}$   
 192 which minimizes the objective function as shown in Eq. (7) [18, 19].

$$193 \quad \sum_{i=1}^n \min_{m=1, \dots, k} d(q_i, o_m) \quad (7)$$

194 where  $n$  is the number of the data points,  $k$  is the number of the clusters,  $q$  is the data point, and  $o$   
 195 is the data point identified as a medoid.

196 PAM consists of two major steps, i.e. build and swap. The first step is to build initial medoids  
 197 by selecting the first medoid as the data point with the minimum sum of the distance to all other  
 198 points and selecting the subsequent medoids by finding the points which minimize Eq. (7). The  
 199 second step repeatedly swap  $i \in \{o_1, o_2, \dots, o_k\}$  with  $j \in \{q_1, q_2, \dots, q_n\}$  if the swap decreases the  
 200 objective significantly until reaching the convergence [18, 19].

201 PAM requires users to provide the number of clusters  $k$  as an input parameter. In the  
 202 proposed strategy, Dunn Index was used to validate the clustering result and determine the  
 203 optimal value of  $k$ . Dunn Index is expressed as the ratio of the smallest inter-cluster distance to  
 204 the largest intra-cluster distance and is defined in Eq. (8) [20].

$$205 \quad D(\phi) = \frac{\min_{C_k, C_l \in \phi, C_k \neq C_l} (\min_{i \in C_k, j \in C_l} d(i, j))}{\max_{C_m \in \phi} (\max_{i, j \in C_m, i \neq j} d(i, j))} \quad (8)$$

206 where  $C_k$  and  $C_l$  are the clusters belong to the set of the identified clusters  $\phi$ . A higher Dunn  
207 Index means a better clustering result. The optimal number of clusters  $k$  was determined based on  
208 the highest Dunn Index within the defined range of the number of clusters.

## 209 2.5 Buildings classification with hierarchical clustering

210 A hierarchical clustering with the heat map visualization technique was used to group  
211 buildings that share the similar daily load characteristics. Hierarchical clustering is a bottom-up  
212 strategy, which starts with placing each object in its own cluster and then merges the atomic  
213 clusters into larger clusters until all objects are in a single cluster [21]. Complete-linkage, which  
214 is the maximum ED of the data objectives in two clusters, was used to measure the distance  
215 between the clusters.

216 An advantage of the hierarchical clustering is that the overall process can be represented by a  
217 tree structure graph called a dendrogram. The dendrogram can help to visualize the cluster  
218 structure and assist in determining the optimal number of clusters. Fig. 3 illustrated a dendrogram  
219 with three data points, where the ordinate axis indicated the distance between the data  
220 points/clusters. The split points indicated the distance between the two data points/clusters. The  
221 higher the split point, the less similarity between the data points/clusters [4]. Clusters can be  
222 determined by the dashed line shown in Fig. 3, which is a user-defined threshold. The data points  
223 under the same split point below the dashed line can be merged into a cluster while the split  
224 points above the dashed line are kept unchanged. For instance, the data points #1 and #3 were  
225 under the same split point and below the dashed line and they will be merged into the same  
226 cluster while the data point #2 formed another cluster. The threshold can be determined  
227 graphically or based on the cluster validation index such as Dunn Index. More details of the  
228 hierarchical clustering can be found in [21].

### 229 3. Performance evaluation of the proposed strategy

230 In this study, the proposed strategy was implemented in R [22] while PAM algorithm was  
231 implemented using the R package cluster [23]. The majority of the figures presented in this study  
232 were generated using R package ggplot2 [24].

#### 233 3.1 Description of the case study buildings

234 The performance of the proposed strategy was evaluated based on the heating energy usage  
235 data collected from 19 higher education buildings, with a total floor area of approximately  
236 200,000 m<sup>2</sup>, at Norwegian University of Science and Technology in Trondheim, Norway. The  
237 hourly building operational data were collected through a web-based Energy Monitoring System.

238 Most of these 19 buildings were built before the year 2000, and the buildings built between  
239 1960 and 1970 accounted for a large part. The energy certificate of the buildings indicated that  
240 the U-values of the exterior walls of the majority buildings were in the range of 0.4-0.60 W/m<sup>2</sup>K,  
241 which failed to comply with the current energy efficiency regulations. Table 1 summarizes the  
242 major information of the 19 buildings studied. More information on these buildings can be found  
243 in [25].

244 The heating demand of these higher education buildings was supplied through a district  
245 heating network and each individual building was equipped with a dedicated heating energy  
246 usage meter. The three-year hourly heating energy data collected from September to April in  
247 2011-2013 were used in this study for performance evaluation of the proposed strategy.

#### 248 3.2 Data pre-processing

249 The generalized ESD test method was first used to detect and remove outliers. Fig. 4  
250 illustrates the three-year hourly heating energy usage data collected from building 03 with the  
251 outliers identified (i.e. red circles). It can be seen that there is a large variation in the heating

252 demand annually. The highest heating demand generally occurred in January and February. It  
253 should be noted that there was a small heating demand from May to August but this amount of  
254 heating demand was significantly lower than that during the main heating period and was  
255 therefore not considered in this study.

256 The data were then standardized to zero mean and one standard deviation and transformed to  
257 daily segments. After removing the daily load profiles with small variations and daily load  
258 profiles in the weekends, a total of 9,062 daily heating energy usage profiles were generated after  
259 the completion of the data pre-processing step.

### 260 3.3 Identification of typical daily heating energy usage profiles

261 The number of clusters selected will directly influence the identification of the typical daily  
262 load profiles. A too small cluster number might result in meaningless typical daily load profiles  
263 while a large cluster number requires a large computational cost and increases the difficulties in  
264 the results evaluation and interpretation. In this study, the optimal cluster number  $k$  (i.e. the  
265 number of the typical daily load profiles) was selected between 5 and 15. Fig. 5 presents Dunn  
266 Index calculated when using different numbers of the clusters. It is shown that the highest Dunn  
267 Index resulted when the cluster number was 11, which was therefore determined as the optimal  
268 cluster number in this study.

269 The boxplot of the aggregated dissimilarity measure of the identified clusters is illustrated in  
270 Fig. 6 for visualization and removal of the daily load profiles beyond the threshold. It can be  
271 observed that the number of the daily load profiles in all clusters ranged from 474 to 1413,  
272 indicating that there was no cluster formed with few daily load profiles. It can also be seen that  
273 all clusters contained the extreme daily load profiles (i.e. black dots) with the aggregated  
274 dissimilarity beyond the threshold (i.e.  $Q3+1.5IQR$ ) and these extreme daily load profiles were

275 removed in subsequent analysis. A total of 8,521 daily load profiles remained after removing the  
276 identified outlier (i.e. extreme daily load profiles) from the dataset. The removal of this small  
277 fraction of the extreme daily load profiles could enhance the visualization of the identified typical  
278 daily load profiles without significant loss of the information.

279 Fig. 7 shows the identified typical daily load profiles by averaging all daily load profiles in  
280 each cluster after the removal of the extreme daily load profiles. The red curves in the figure  
281 showed the typical daily load profiles identified while the gray curves were all corresponding  
282 daily load profiles in this cluster. It can be found that there was a clear boundary in the heating  
283 demand between the working hours and non-working hours in some typical daily load profiles  
284 such as the load profiles 2 and 8 while that in some typical daily load profiles (e.g. the load  
285 profiles 1 and 5) were not very clear. There was no obvious boundary in the load profile 10.  
286 Moreover, the nighttime from 22:00 to 03:00 of next day was the lowest heating energy usage  
287 period for the majority of the typical daily load profiles identified except the typical load profiles  
288 6, 7 and 10 with a noticeable high heating demand during the nighttime which is worthwhile for  
289 further investigation.

290 Fig. 8 shows the weekday distribution of the building daily load profiles in the identified  
291 clusters, in which y-axis represents the percentage of the number of days belongs to each  
292 weekday to the total number of days in each cluster. It was shown that the daily load profiles on  
293 Tuesday, Wednesday, Thursday and Friday in each cluster were almost evenly distributed. In  
294 some clusters such as the clusters 4 and 11, the number of days on Monday was obviously  
295 different from that on the other weekdays and the reason behind this is presented in Section 4.  
296 Therefore, this weekday load profile distribution can assist in determining whether a specific load  
297 profile existed only in some specific days of a week.

298 Table 2 summarizes the key characteristics and the estimated high heating demand period of  
299 the typical daily load profiles identified. To understand the knowledge and information  
300 discovered by the proposed strategy, the profiles with a relatively high demand in the early  
301 morning and late night as well as those with clear heating demand peaks and troughs will be  
302 further investigated in Section 4. These include the typical daily load profiles 4, 6, 7, 9 and 11.  
303 The rest of the typical load profiles were either similar to the typical daily load profiles  
304 mentioned above or did not contain interesting characteristics and were therefore not further  
305 investigated in this study.

#### 306 3.4 Building classification based on the identified typical daily load profiles

307 In this section, 19 case study buildings were grouped according to the typical daily heating  
308 load profiles identified. In order to eliminate the influence from the insignificant profiles, the first  
309 two most dominant profiles of each building (see Table 3) were selected as the features for  
310 building classification. From Table 3, it can be seen that for some buildings such as buildings 02,  
311 14 and 17, the most dominant profile accounted for a large proportion of the total number of days  
312 remained for the typical daily load profile identification. For instance, 436 days out of 490 of  
313 building 02 were in the most dominant profile, demonstrating that the daily load variation of this  
314 building was consistent. In contrast, the number of days in the most dominant profiles of some  
315 buildings such as building 10 and 15 were relatively small, which indicated that these buildings  
316 did not have a consistent daily load profile during the time period investigated (2011-2013).

317 The percentages of the first two most dominant profiles were then used to group the buildings  
318 that share the same daily energy usage characteristics based on the hierarchical clustering. Fig. 9  
319 presents the dendrogram of building classification results, in which the buildings in the same  
320 cluster were marked with the same color. In this study, the threshold (i.e. dashed line in the figure)

321 was visually selected due to the small number of the data points (i.e. buildings) used. It can be  
322 seen that some clusters were formed with a single building while some clusters were formed with  
323 several buildings. For instance, building 02 was identified as an individual cluster and buildings  
324 01, 05, 10, 11 and 12 were grouped into one single cluster.

325 In order to better visualize and confirm the clustering results, the number of days belongs to a  
326 typical daily load profile of different buildings were plotted as a heat map and are shown in Fig.  
327 10. In this figure, the relative proportion ( $RP$ ) was determined using Eq. (9) and the same order of  
328 the building number as illustrated in Fig. 9 was used. It was visually shown that the majority of  
329 the buildings had one significant dominant profile.

$$330 \quad RP = \frac{N_d}{N_{d,max}} \quad (9)$$

331 where  $N_d$  stands for the number of days belongs to a typical daily load profile of an individual  
332 building and  $N_{d,max}$  stands for the maximum number of days belong to a typical daily load profile  
333 of the same building.

#### 334 **4. Interpretation of the identified typical daily load profiles**

335 In order to understand the reasons behind the main characteristics of the typical daily load  
336 profile identified, buildings 02, 14, 17, 08 and 03 were selected based on the clustering results  
337 and used to represent the typical daily load profiles of 4, 6, 7, 9 and 11 presented in Fig. 7,  
338 respectively.

##### 339 **4.1 Building 02 – Typical daily load profile 4**

340 Building 02 is an office and laboratory building which was built in 1965. A recent survey  
341 indicated that this building was poorly insulated with a U-value of 0.91 W/m<sup>2</sup>K for the exterior  
342 wall insulation and a U-value of 0.59 W/m<sup>2</sup>K for the roof insulation. Different from many other  
343 buildings using hot water radiators for space heating, the heating of this building was supplied

344 through ventilation without using heat recovery. However, the heat recovery has been  
345 mandatorily required for decades in Norway in ventilation.

346 Fig. 11(a) shows the heating energy usage of building 02 in the two consecutive days. It was  
347 clearly shown that the high heating demand started at 04:00 in the morning, which was consistent  
348 with the typical daily load profile 4. However, it was much earlier than the normal building  
349 occupied hours. The feedback from the building operator indicated that the occupants in this  
350 building continuously complained about the thermal comfort during the morning time. The  
351 heating period was therefore extended in order to satisfy the occupant thermal comfort and to  
352 provide freezing protection [26].

#### 353 4.2 Building 14 – Typical daily load profile 6

354 Building 14 is a sports center, which was usually operated till midnight. The heating demand  
355 of this building in the two consecutive days is illustrated in Fig. 11(b). The major characteristics  
356 of the two-day heating demand matched well with that of the typical daily load profile 6. The  
357 highest heating demand generally occurred around 19:00. This high heating demand was  
358 probably related to the hot water usage for the shower requirement. The water usage data of this  
359 building in the same two days are presented in Fig. 12. It was clearly shown that there was a high  
360 peak of the water usage at around 19:00, which was in line with the heating energy usage profiles.  
361 It was also found that the water usage of this building dropped to zero at 01:00 which also  
362 matched with the heating demand variation.

#### 363 4.3 Building 17 – Typical daily load profile 7

364 Building 17 is a multi-functional building with offices, educational rooms and laboratories,  
365 which was constructed around the year 1996. As shown in Fig. 11(c), the two-day heating load  
366 profile of this building was similar to that of the typical daily load profile 7 identified. The high  
367 heating demand period lasted till to 23:00. The feedback from the building operator indicated that



368 the building occupants required the building to be heated till to 23:00 for special activity  
369 requirements.

#### 370 4.4 Building 08 – Typical daily load profile 9

371 Building 08 is an old building constructed in 1924 and is also a multi-functional building  
372 with offices, educational rooms, and laboratories. A clear peak and a clear trough can be  
373 observed in Fig. 11(d) at 05:00 and 21:00 respectively, which were consistent with the  
374 information presented in the typical daily load profile 9. The heating demand peak and trough  
375 were found to be mainly caused by the sudden change of the supply water temperature. The  
376 recorded data showed that the hot water was supplied at about 70°C during the daytime and 40°C  
377 during the nighttime. The sudden rise of the supply water temperature in the early morning  
378 resulted in the heating demand peak of the building while the sudden drop of the supply water  
379 temperature in the nighttime led to the occurrence of the trough in the heating load profile. This  
380 relationship between the heating energy usage and the variation in the supply water temperature  
381 was also observed in a previous study [27].

382 The building operator was also approached for the reason why the high heating demand  
383 started at around 05:00. However, no information on this was recorded. This is probably also due  
384 to the poor insulation of the building (i.e. U-value of 1.0 W/m<sup>2</sup>K for the exterior wall insulation  
385 and U-value of 0.7 W/m<sup>2</sup>K for the roof insulation), which might result in a longer pre-heating  
386 period before the building was occupied.

#### 387 4.5 Building 03 – Typical daily load profile 11

388 Building 03 is a mix of offices and laboratories, which was constructed in 1951. The typical  
389 daily load profile 11 was very similar to the typical daily load profile 9. However, in the typical  
390 daily load profile identified, there were very few days from Monday. Fig. 11(e) illustrates the  
391 heating demand of building 03 in two days of Monday and Tuesday. It was clearly shown that

392 there was a heating demand peak at 07:00 and a trough at 17:00 in the daily heating load profile  
393 on Tuesday, which matched well with the typical daily load profile 11. However, on Monday, the  
394 heating demand peak occurred at 06:00. This is mainly due to the fact that, during the weekend,  
395 the heating system was either not running or running with a lower supply water temperature,  
396 resulting in a lower indoor temperature than during the weekdays. In order to achieve a desirable  
397 thermal comfort on Monday morning, the building was therefore pre-heated earlier than that  
398 during the weekdays.

## 399 **5. Comparison between the use of ED-based and PCC-based clustering**

400 In this section, the results of using the ED-based and PCC-based clustering were compared  
401 and presented. The same data pre-processing used for the PCC-based clustering was performed  
402 for the ED-based clustering while the commonly used k-means and ED-based dissimilarity  
403 measure were used to replace PAM and PCC-based dissimilarity measure. The optimal number  
404 of clusters for the ED-based clustering determined was 10, as shown in Figure 13, which was also  
405 determined based on Dunn index.

406 Based on the optimal number of clusters determined, the typical daily heating load profiles  
407 can then be identified after removal of the extreme daily load profiles based on the box plot  
408 analysis. Fig. 14 presents the clustering results and the identified typical daily heating load  
409 profiles using the ED-based clustering. It can be seen that the profiles identified using the k-  
410 means clustering with ED-based dissimilarity measure can still provide some useful information  
411 in the identified typical daily heating load profiles. For instance, a morning peak was observed  
412 and the building was heated till to midnight in the typical daily load profile 7, which was very  
413 similar to the typical daily load profile 6 identified using the proposed strategy.

414 However, some profiles identified such as the typical daily load profiles 3 and 9 were too flat  
415 and cannot provide useful information for further analysis. Some important information, e.g.

416 05:00 heating demand peak (corresponding to the load profile 9 in Fig. 7), 04:00 high heating  
417 demand start time (corresponding to the load profile 4 in Fig. 7), 17:00 low trough  
418 (corresponding to load profile 11 in Fig. 7), identified by the proposed strategy cannot be  
419 identified using the k-means clustering with ED-based dissimilarity measure. In addition, some  
420 profiles such as the load profiles 2 & 7, and the load profiles 6 & 8 presented in Fig. 14 showed  
421 very similar trends but with different magnitudes. This further demonstrated that the ED-based  
422 dissimilarity measure tends to identify daily load profiles that were similar in terms of the  
423 intensity.

424 The heat map in Fig. 15 illustrated the number of days belongs to a typical daily load profile  
425 of different buildings when using the ED-based clustering. The order of the buildings in the map  
426 was also determined based on the hierarchical clustering. Compared to the results presented in  
427 Fig. 10, it was clearly shown that the building heating energy usage cannot be characterized by  
428 the most dominant load profiles as there was no clear difference between the number of days  
429 belong to the most dominant typical daily load profile and that belongs to the other typical daily  
430 load profiles. It was also demonstrated that it is difficult to use the ED-based clustering identified  
431 typical daily load profiles for building classification.

## 432 **6. Conclusions**

433 Understanding multiple buildings energy performance requires advanced data analytics. This  
434 paper presented a variation focused Partitioning Around Medoids (PAM) cluster analysis strategy  
435 to identify the typical daily load profiles of higher education buildings, in which Pearson  
436 Correlation Coefficient was used as the dissimilarity measure to group the daily load profiles on  
437 the basis of the variation similarity instead of the magnitude similarity.

438 The performance of the proposed strategy was evaluated using the heating energy usage data  
439 of 19 higher education buildings in Norway collected from 2011 to 2013. The results showed that

440 the proposed strategy can identify and discover the information related to building daily heating  
441 energy usage characteristics, including daily high heating demand start time and end time, the  
442 peaks, troughs and variations of daily heating energy usage. The results obtained also confirmed  
443 the effectiveness of the proposed strategy in identifying the typical daily heating energy usage  
444 profiles in terms of the variation similarity.

445 The identified daily heating energy usage characteristics can be used to assist in the  
446 development of advanced building control and fault detection & diagnosis strategies, and cost-  
447 effective demand side management techniques. The information discovered is also useful to  
448 support the energy planning and retrofitting of higher education buildings. This method could be  
449 adapted to identify the daily energy usage characteristics of other types of buildings.

450

## 451 **Acknowledgement**

452 This research work was made possible through an Endeavour Research Fellowship. The first  
453 author would like to thank the support of Australian Government - Department of Education and  
454 Training.

## 455 **References**

- 456 [1] Z. Ma, P. Cooper, D. Daly, L. Ledo. Existing building retrofits: Methodology and state-of-  
457 the-art. *Energy and buildings*. 55 (2012) 889-902.
- 458 [2] L. Pérez-Lombard, J. Ortiz, C. Pout. A review on buildings energy consumption information.  
459 *Energy and Buildings*. 40 (2008) 394-8.
- 460 [3] H. Gadd, S. Werner. Heat load patterns in district heating substations. *Applied Energy*. 108  
461 (2013) 176-83.
- 462 [4] P.R.S. Jota, V.R.B. Silva, F.G. Jota. Building load management using cluster and statistical  
463 analyses. *International Journal of Electrical Power & Energy Systems*. 33 (2011) 1498-505.

- 464 [5] C.M.R. do Carmo, T.H. Christensen. Cluster analysis of residential heat load profiles and the  
465 role of technical and household characteristics. *Energy and Buildings*. 125 (2016) 171-80.
- 466 [6] A.K. Jain. Data clustering: 50 years beyond K-means. *Pattern Recognition Letters*. 31 (2010)  
467 651-66.
- 468 [7] C. Flath, D. Nicolay, T. Conte, C. van Dinther, L. Filipova-Neumann. Cluster analysis of  
469 smart metering data: An implementation in practice. *Business & Information Systems  
470 Engineering*. 4 (2012) 31-9.
- 471 [8] C. Miller, Z. Nagy, A. Schlueter. Automated daily pattern filtering of measured building  
472 performance data. *Automation in Construction*. 49 (2015) 1-17.
- 473 [9] M.P. Fernandes, J.L. Viegas, S.M. Vieira, J.M. Sousa. Analysis of residential natural gas  
474 consumers using fuzzy c-means clustering. 2016 IEEE International Conference on Fuzzy  
475 Systems (FUZZ). pp. 1484-91.
- 476 [10] I.P. Panapakidis, T.A. Papadopoulos, G.C. Christoforidis, G.K. Papagiannis. Pattern  
477 recognition algorithms for electricity load curve analysis of buildings. *Energy and Buildings*. 73  
478 (2014) 137-45.
- 479 [11] L. Kaufman, P.J. Rousseeuw. *Finding groups in data: An introduction to cluster analysis*.  
480 John Wiley & Sons, Inc., Hoboken, New Jersey, USA, 2008.
- 481 [12] L. Ledo. *Energy efficiency and thermal comfort upgrades for higher education buildings*.  
482 University of Wollongong, PhD thesis, 2015.
- 483 [13] P.N. Tan, M. Steinbach, V. Kumar. *Introduction to data mining*. Pearson Addison Wesley,  
484 USA, 2006.
- 485 [14] J.E. Seem. Using intelligent data analysis to detect abnormal energy consumption in  
486 buildings. *Energy and Buildings*. 39 (2007) 52-8.

- 487 [15] C. Fan, F. Xiao, S. Wang. Development of prediction models for next-day building energy  
488 consumption and peak power demand using data mining techniques. *Applied Energy*. 127 (2014)  
489 1-10.
- 490 [16] I. Khan, A. Capozzoli, S.P. Corgnati, T. Cerquitelli. Fault detection analysis of building  
491 energy consumption using data mining techniques. *Energy Procedia*. 42 (2013) 557-66.
- 492 [17] H. Blockeel, L. De Raedt, J. Ramon. Top-down induction of clustering trees. *Proceedings of*  
493 *the Fifteenth International Conference on Machine Learning*, 1998. pp. 55-63.
- 494 [18] L. Kaufman, P. Rousseeuw. *Clustering by means of medoids*. The First International  
495 *Conference on Statistical Data Analysis based on the L<sub>1</sub>-Norm and Related Methods*, Neuchatel,  
496 Switzerland, 1987.
- 497 [19] A. Struyf, M. Hubert, P. Rousseeuw. Clustering in an object-oriented environment. *Journal*  
498 *of Statistical Software*. 1 (1997) 1-30.
- 499 [20] G. Brock, V. Pihur, S. Datta, S. Datta. *clValid: An R package for cluster validation*. *Journal*  
500 *of Statistical Software*. 25 (2008) 1-22.
- 501 [21] J. Han, M. Kamber, J. Pei. *Data mining: Concepts and techniques (Second Edition)*. Morgan  
502 Kaufmann Publishers, USA, 2006.
- 503 [22] R Development Core Team. *R: A language and environment for statistical computing*. R  
504 *Foundation for Statistical Computing*, Vienna, Austria, 2008.
- 505 [23] M. Maechler, P. Rousseeuw, A. Struyf, M. Hubert, K. Hornik. *cluster: Cluster analysis*  
506 *basics and extensions*. 2015.
- 507 [24] H. Wickham. *ggplot2: Elegant Graphics for Data Analysis*. Springer, New York, 2009.
- 508 [25] J. Guan, N. Nord, S. Chen. Energy planning of university campus building complex: Energy  
509 usage and coincidental analysis of individual buildings with a case study. *Energy and Buildings*.  
510 124 (2016) 99-111.

511 [26] N. Nord, M.H. Martin. Freeze protection method in ventilation system using two hydronic  
512 circuits. The 12<sup>th</sup> International Conference on Air Distribution in Rooms, Trondheim, Norway,  
513 2011.

514 [27] N. Djuric, V. Novakovic, F. Frydenlund. Heating system performance estimation using  
515 optimization tool and BEMS data. Energy and Buildings. 40 (2008) 1367-76.

516

517

Table 1 Major information of the case study buildings

Building NO.	Construction year	Main functions <sup>#</sup>	Floor area (m <sup>2</sup> )	Building NO.	Construction year	Main function	Floor area (m <sup>2</sup> )
01	1962	O/E/L	15,026	11	1968	O/L	12,861
02	1965	O/L	3,030	12	1910	O	3,375
03	1951	O/L	2,215	13	1981	O/E/L	3,955
04	1960	O/E/L	7,598	14	1966	S	4,046
05	1966	O/E/L	11,400	15	1975	O/E/L	18,175
06	1958	O/E/L	12,600	16	1951	O/E/L	5,053
07	1965	O/E/L	9,168	17	1996	O/E/L	2,476
08	1924	O/E/L	4,116	18	2002	E/L	4,312
09	1960	O/L	5,028	19	2000	O/E/L	52,773
10	1961	O/L	17,936				

519 # O: office; E: educational room; L: laboratory; S: sports complex.

520 Table 2 Key characteristics of the identified typical daily heating energy usage profiles

Typical load profile No.	Est. high heating demand period	Weekday load profile almost evenly distributed	Main characteristics
1	07:00-15:00	Yes	There was a high heating demand from 07:00 to 10:00. The heating demand was then gradually decreased till to 16.00 and then kept relatively stable.
2	07:00-17:00	No	The high heating demand occurred during the office hours. A clear heating demand peak can be observed at 07:00.
3	07:00-18:00	No	There was a clear heating demand peak at 07:00 and the heating demand was then gradually decreased till to 18:00.
4	04:00-17:00	No	A high heating demand started at around 04:00 and then kept relatively stable till to 17:00.
5	07:00-18:00	Yes	The daily heating demand variations were similar to that of the typical load profile 1.
6	09:00-24:00	Yes	There was a small peak at 06:00. A high heating demand started at 09:00 and lasted till to the midnight.
7	09:00-23:00	Yes	Similar to the load profile 6 but the heating demand during the high heating demand period was more stable.
8	06:00-18:00	No	Similar to the load profile 2. However, there was a clear trough at 19:00.
9	05:00-20:00	No	Similar to the load profiles 2 and 8 but there was a clear peak at 05:00 and a clear trough at 21:00.
10	Not clear	Yes	The heating demand during 24 hours was relatively stable. However, the demand in the early morning was slightly higher than the rest of the day.
11	07:00-16:00	No	Similar to the load profiles 2, 8 and 9. There was a clear heating demand peak at 07:00 and a clear trough at 17:00.



521

Table 3 Summary of the first two most dominant profiles of individual buildings

Building No.	Total number of days	The most dominant profile			The 2 <sup>nd</sup> most dominant profile		
		Typical daily load profile No.	Total days	Percentage (%)	Typical daily load profile No.	Total days	Percentage (%)
1	457	5	240	53	3	79	17
2	490	4	436	89	3	17	3
3	486	11	252	52	8	57	12
4	458	11	207	45	4	58	13
5	436	5	271	62	7	36	8
6	437	3	246	56	1	93	21
7	471	8	300	64	11	65	14
8	471	9	285	61	7	87	18
9	448	10	124	28	3	118	26
10	439	7	107	24	5	103	23
11	471	5	175	37	7	81	17
12	371	5	172	46	10	83	22
13	449	3	148	33	2	127	28
14	495	6	440	89	7	35	7
15	382	1	94	25	3	73	19
16	486	2	316	65	4	67	14
17	480	7	386	80	5	40	8
18	367	7	152	41	6	115	31
19	427	1	116	27	5	112	26

522

523

524

525 **Figure Captions**

526 Fig. 1 Outline of the variation focused cluster analysis strategy.

527 Fig. 2 Comparison between the PCC and ED-based dissimilarity measures (a) ED; (b)&(c) PCC.

528 Fig. 3 Illustration of the dendrogram with three data points.

529 Fig. 4 Illustration of the building heating energy usage and outliers identified - building 03.

530 Fig. 5 Dunn Index calculated for different numbers of the clusters - PCC-based clustering.

531 Fig. 6 Boxplot of the aggregated dissimilarities of the identified clusters.

532 Fig. 7 Typical daily heating load profiles (red) identified using the proposed strategy with all  
533 corresponding daily load profiles (gray).

534 Fig. 8 Weekday load profile distribution in different clusters identified.

535 Fig. 9 Dendrogram of building classification results.

536 Fig. 10 Heat map of the typical daily load profiles in different buildings - PCC-based clustering.

537 Fig. 11 Illustrations of the heating energy usage of the buildings in two consecutive days.

538 Fig. 12 Water usage of building 14.

539 Fig. 13 Dunn Index calculated for different numbers of the clusters - ED-based clustering.

540 Fig. 14 Typical daily heating load profiles (red) identified using the ED-based clustering with all  
541 corresponding daily load profiles (gray).

542 Fig. 15 Heat map of the typical daily load profiles in different buildings - ED-based clustering.

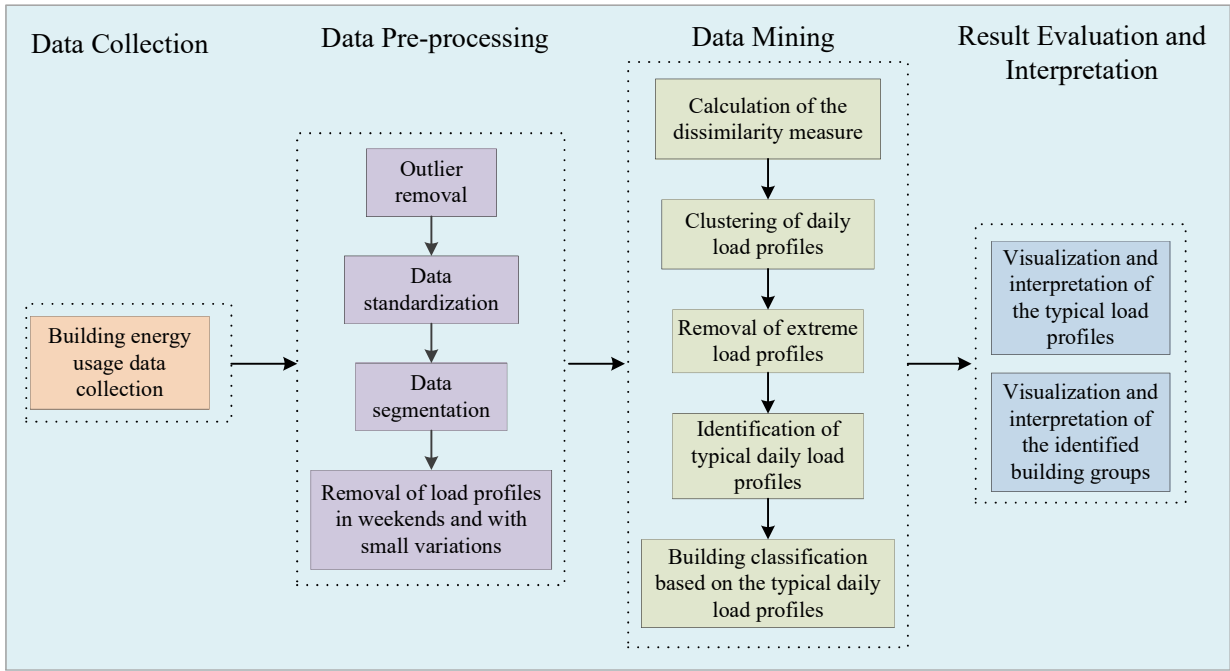
543

544

545

546

547

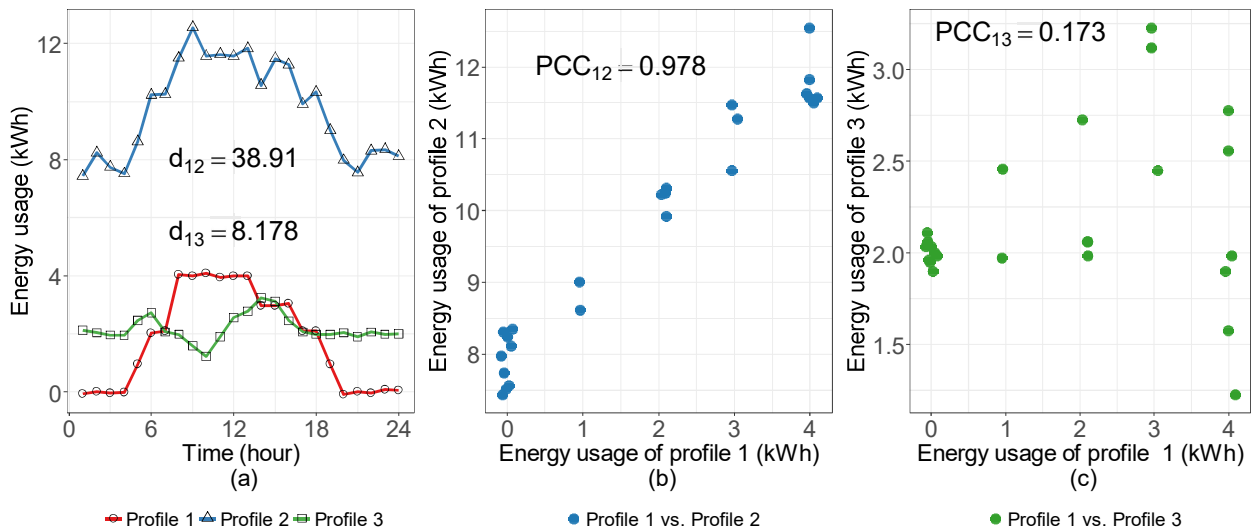


548

549

Fig. 1 Outline of the variation focused cluster analysis strategy.

550

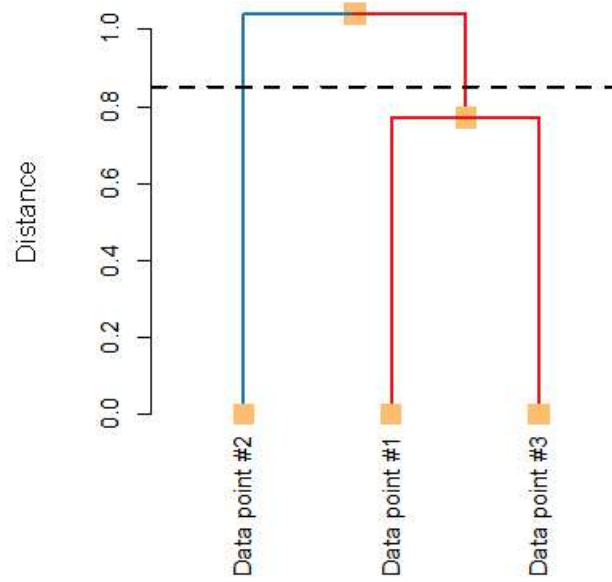


551

552

Fig. 2 Comparison between the PCC and ED-based dissimilarity measures (a) ED; (b)&(c) PCC.

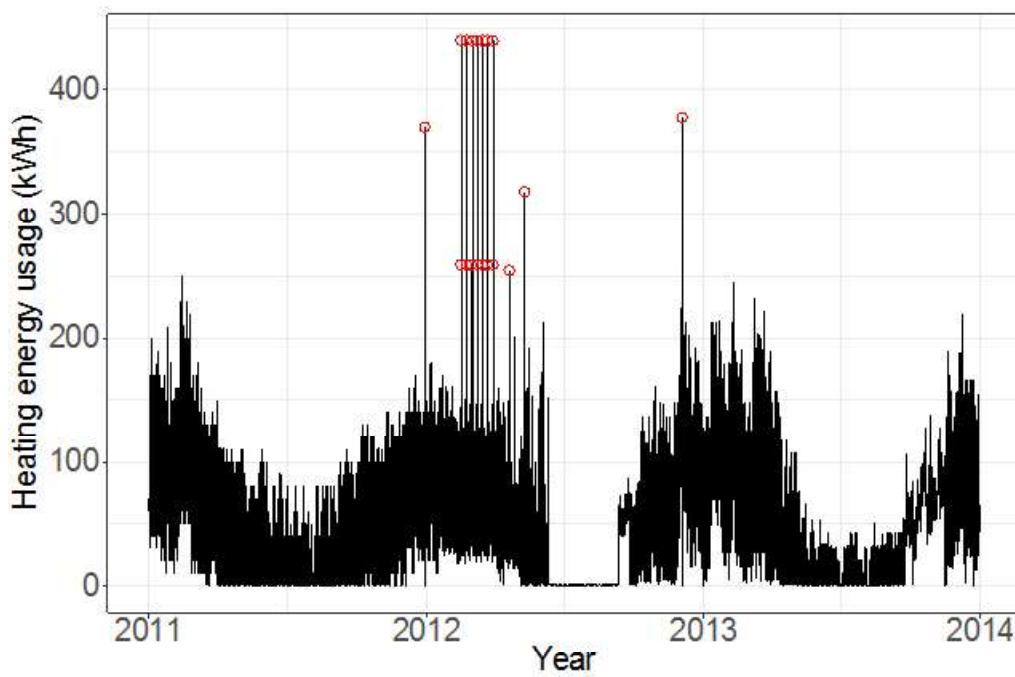
553



554

555

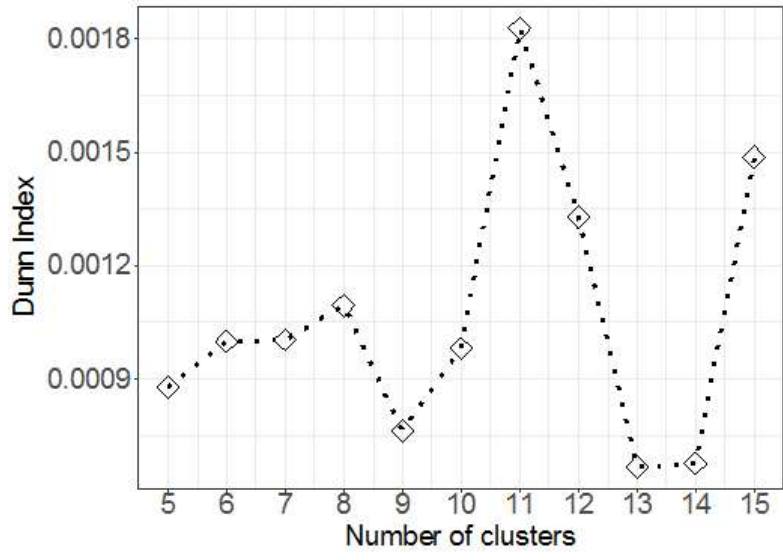
Fig. 3 Illustration of the dendrogram with three data points.



556

557

Fig. 4 Illustration of the building heating energy usage and outliers identified - building 03.



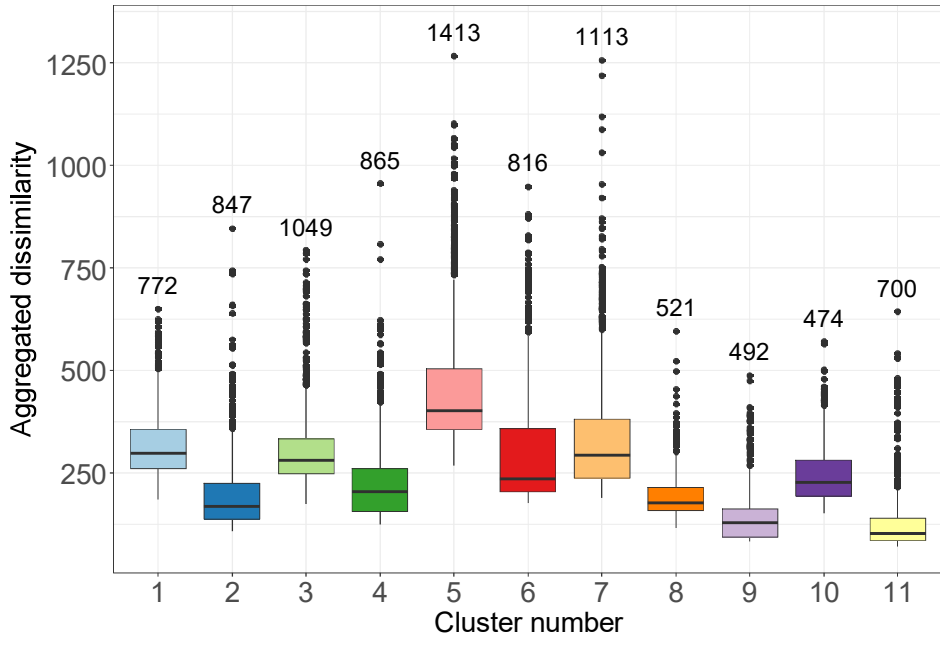
558

559

Fig. 5 Dunn Index calculated for different numbers of the clusters - PCC-based clustering.

560

561

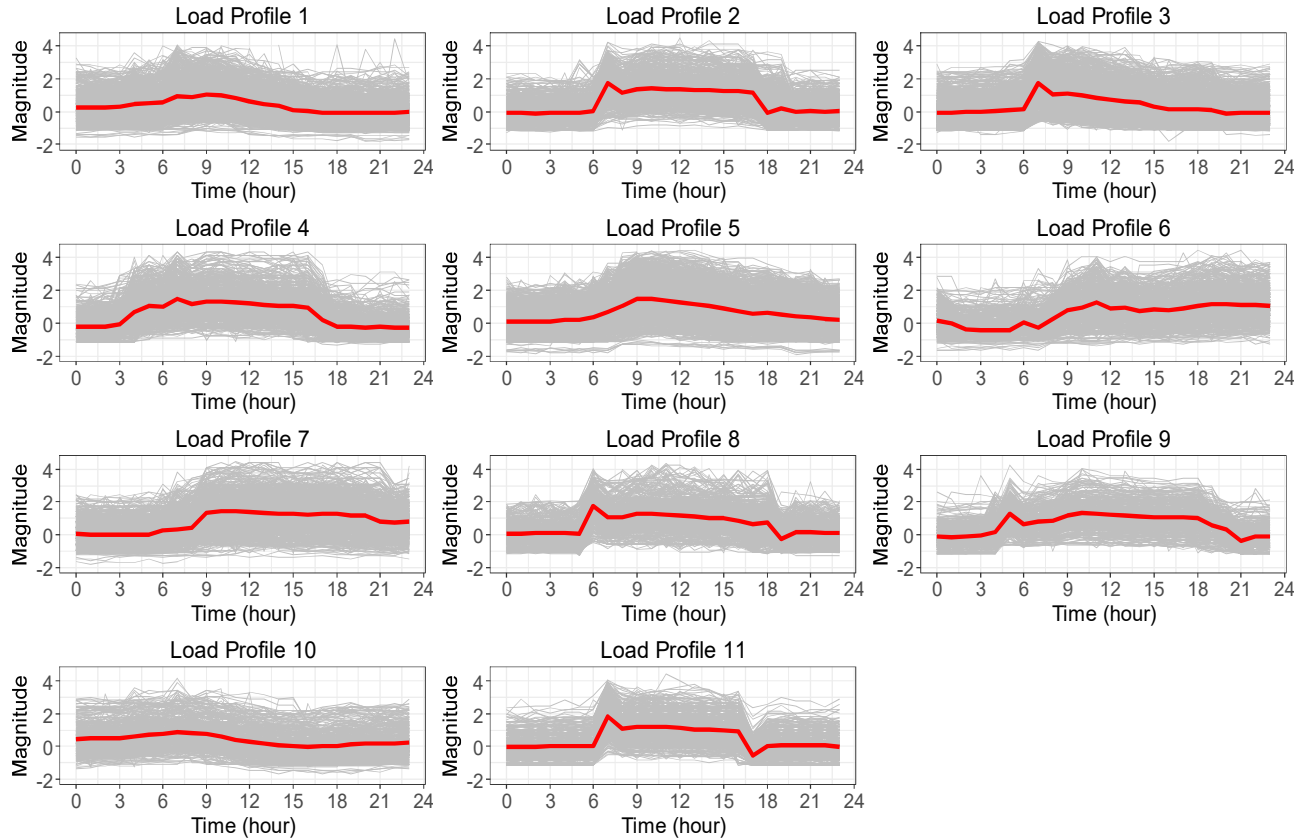


562

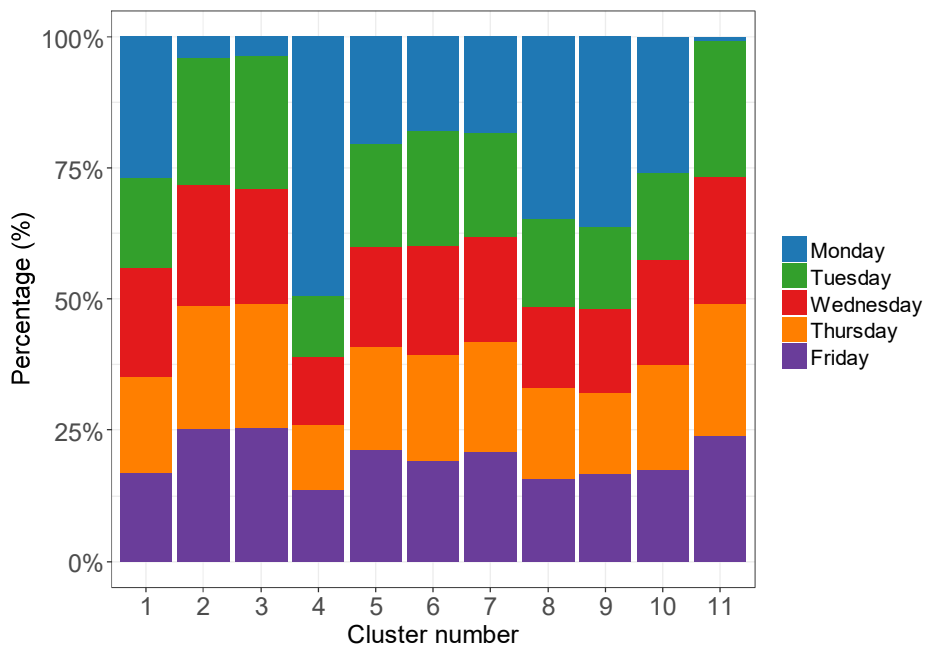
563

Fig. 6 Boxplot of the aggregated dissimilarities of the identified clusters.

564

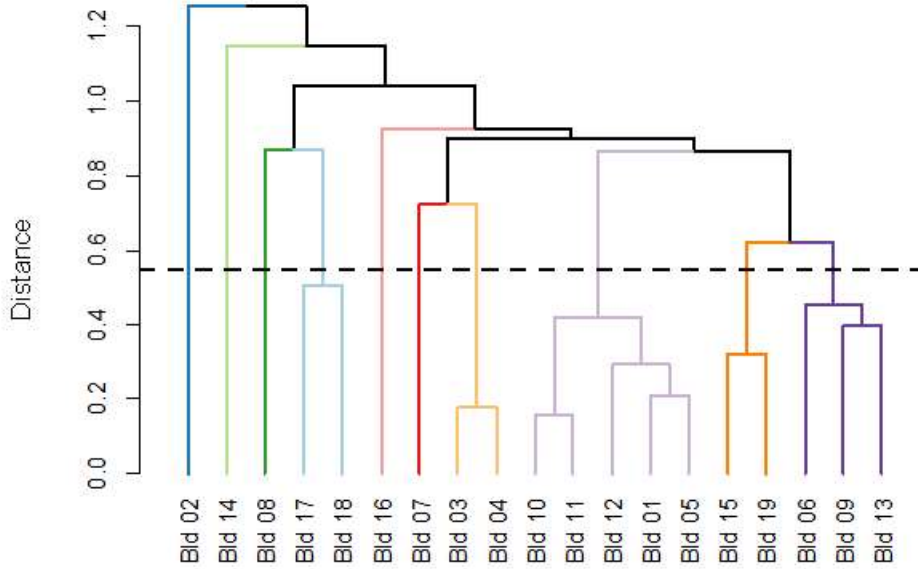


565  
 566 Fig. 7 Typical daily heating load profiles (red) identified using the proposed strategy with all  
 567 corresponding daily load profiles (gray).  
 568



569  
 570 Fig. 8 Weekday load profile distribution in different clusters identified.

571

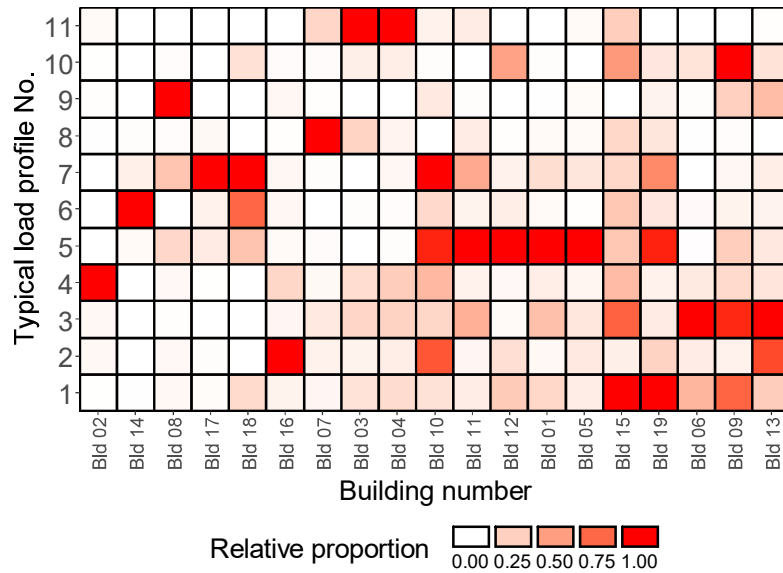


572

573

Fig. 9 Dendrogram of building classification results.

574



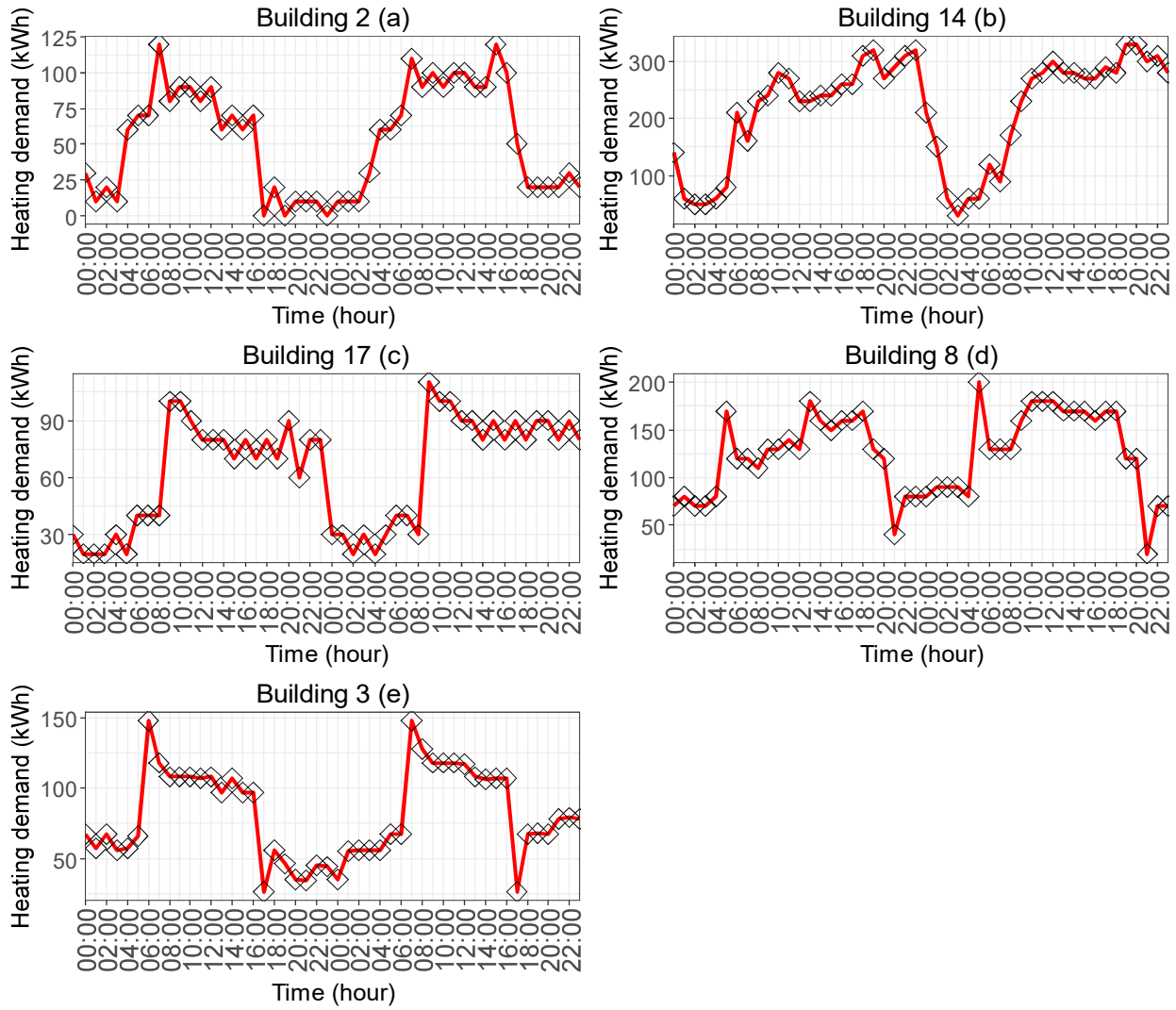
575

576

Fig. 10 Heat map of the typical daily load profiles in different buildings - PCC-based clustering.

577

578



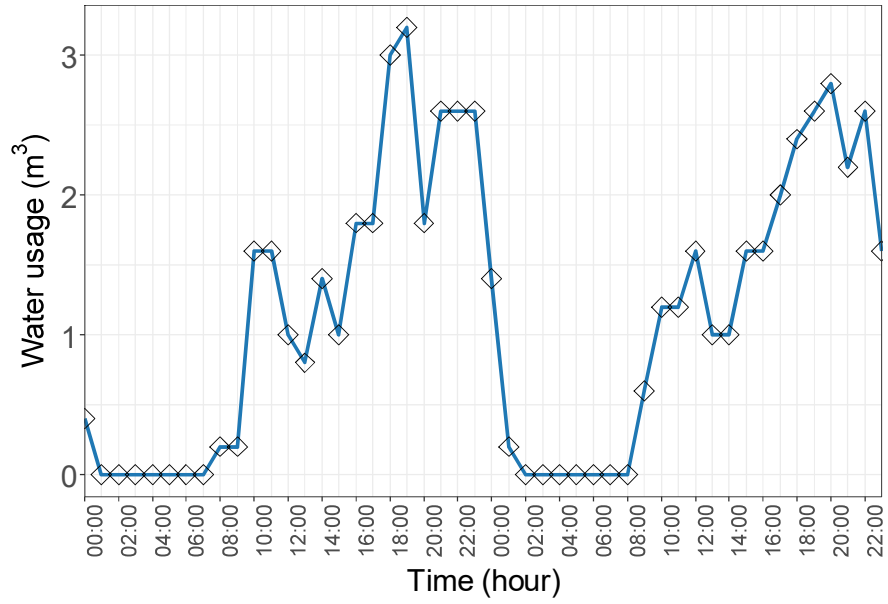
579

580

581

Fig. 11 Illustrations of the heating energy usage of the buildings in two consecutive days.

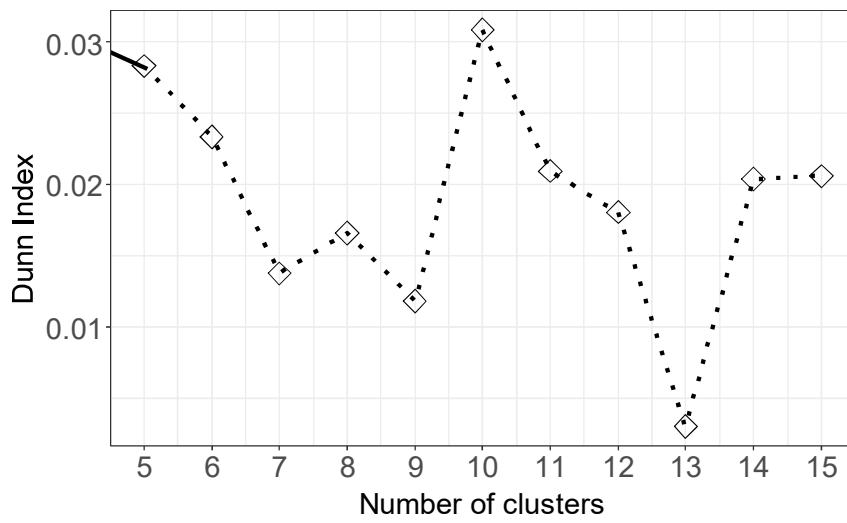




582

583

Fig. 12 Water usage of building 14.



584

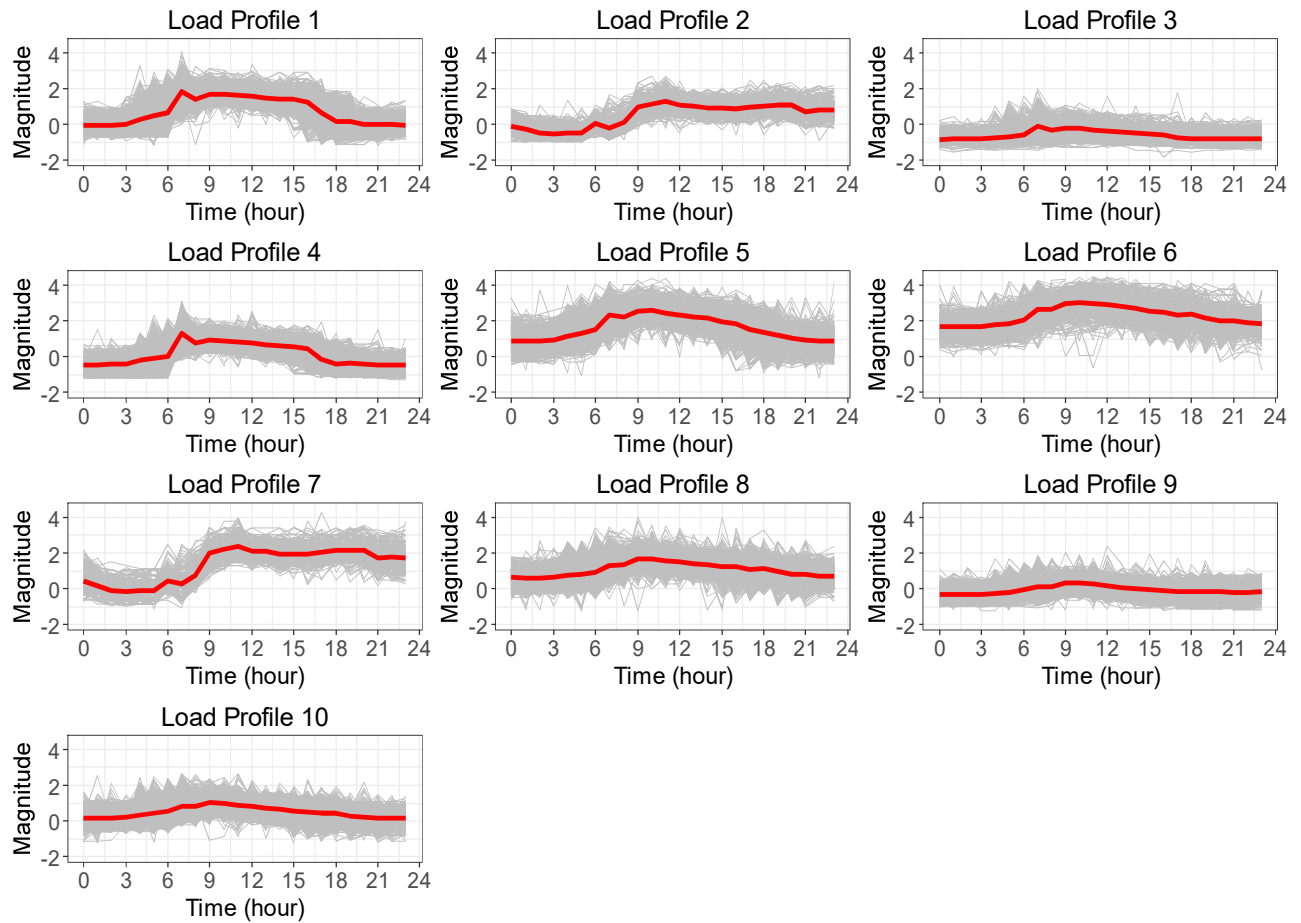
585

Fig. 13 Dunn Index calculated for different numbers of the clusters - ED-based clustering.

586

587

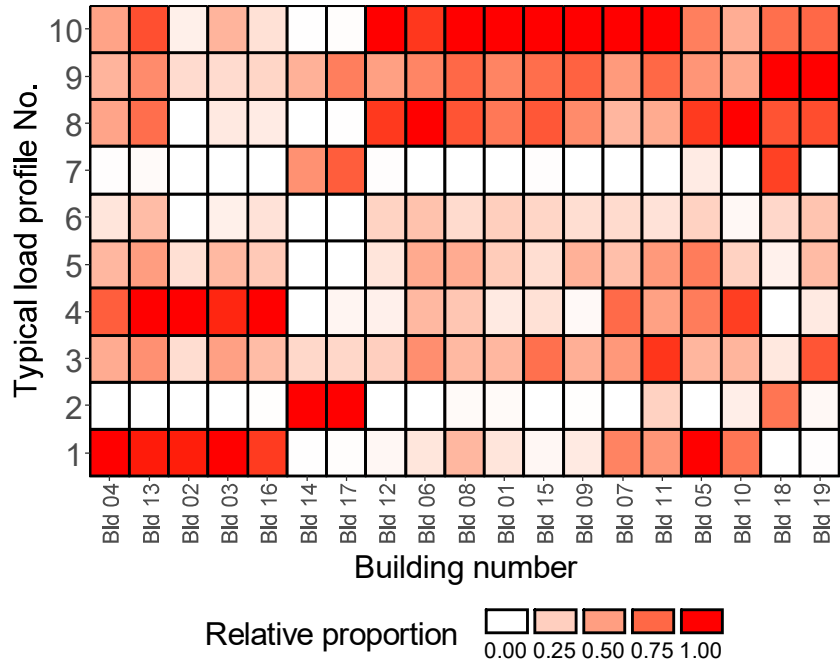
588



589

590 Fig. 14 Typical daily heating load profiles (red) identified using the ED-based clustering with all  
 591 corresponding daily load profiles (gray).

592



593  
594

Fig. 15 Heat map of the typical daily load profiles in different buildings - ED-based clustering.