

Journal Pre-proofs

A reduced-dimension feature extraction method to represent retail store electricity profiles

Ramon Granell, Colin J. Axon, Maria Kolokotroni, David C.H. Wallom

PII: S0378-7788(22)00679-X

DOI: <https://doi.org/10.1016/j.enbuild.2022.112508>

Reference: ENB 112508

To appear in: *Energy & Buildings*

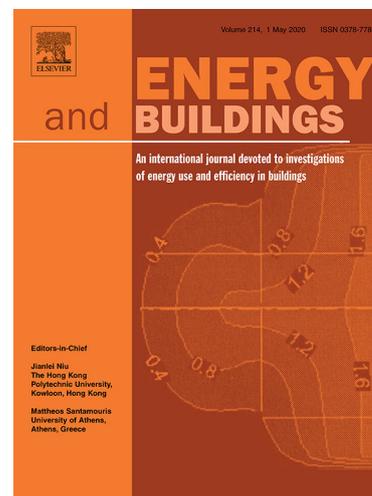
Received Date: 19 May 2022

Accepted Date: 20 September 2022

Please cite this article as: R. Granell, C.J. Axon, M. Kolokotroni, D.C.H. Wallom, A reduced-dimension feature extraction method to represent retail store electricity profiles, *Energy & Buildings* (2022), doi: <https://doi.org/10.1016/j.enbuild.2022.112508>

This is a PDF file of an article that has undergone enhancements after acceptance, such as the addition of a cover page and metadata, and formatting for readability, but it is not yet the definitive version of record. This version will undergo additional copyediting, typesetting and review before it is published in its final form, but we are providing this version to give early visibility of the article. Please note that, during the production process, errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.

© 2022 Published by Elsevier B.V.



13 A reduced-dimension feature extraction method to
 14 represent retail store electricity profiles

15 Ramon Granell^{a,b,1}, Colin J. Axon^a, Maria Kolokotroni^a, David C.H.
 16 Wallom^b

17 ^a*Institute of Energy Futures, Brunel University London, Uxbridge, London UB8 3PH,*
 18 *United Kingdom.*

19 ^b*Oxford e-Research Centre, Department of Engineering Science, University of Oxford, 7*
 20 *Keble Road, Oxford OX1 3QG, United Kingdom.*

21 **Abstract**

22 Characterising the inter-seasonal energy performance of buildings is a
 23 useful tool for a business to understand what is normal for its portfolio of
 24 premises and to detect anomalous patterns of energy demand. When adding
 25 a new building to the portfolio, it will be useful to predict what will be
 26 the likely energy use as part of on-going monitoring of the site. For a large
 27 portfolio of buildings with, say, half-hourly energy use measurements (48 di-
 28 mensions), analysis and prediction will require machine learning tools. Even
 29 so, it is advantageous to minimise the amount of data and number of di-
 30 mensions and features required to find useful patterns in the measurement
 31 stream. Our aim is to devise a reduced feature set that can generate a sta-
 32 tistically reasonable representation of daily electricity load profiles of retail
 33 stores and small supermarkets. We then test if our method is sufficiently
 34 accurate to predict and cluster measured patterns of demand. We propose
 35 an automatic method to extract features such as times and average demands
 36 from electricity load profiles. We used four regression models for prediction
 37 and six clustering methods to compare with the results obtained using all of
 38 the readings in the load profile. We found that the reduced feature set gave a
 39 good representation of the load profile, with only small prediction and clus-
 40 tering errors. The results are robust as prediction is supervised learning and
 41 clustering is unsupervised. This simplified feature set is a concise way to rep-
 42 resent profiles without using small variances of the demand that do not add
 43 useful information to the overall picture. As modern sensor systems increase
 44 the volume, availability, and immediacy of data, using reduced dimensional

45 datasets will be key to extracting useful information from high-resolution
46 data streams.

47 *Keywords:* clustering, electricity demand, commercial, prediction, machine
48 learning, supermarket

49 **Abbreviations**

50 **ANN** Artificial neural networks

51 **ED** Euclidean distance

52 **EDLP** Electricity daily load profile

53 **kNNR** k-nearest neighbours regression

54 **ML** Machine learning

55 **NP** Normalised percentage difference with respect to the original EDLP

56 **OLS** Ordinary least of squares

57 **SE** Supermarkets using only electricity

58 **SEG** Supermarkets using electricity and gas

59 **SVR** Support vector regression

60 **Symbols**

61 e_i electricity consumed (kWh) between the $(i - 1)$ -th and i -th time interval

62 k number of EDLPs used for the prediction

63 p number of previous years used to predict the EDLP

64 s_0 off-peak time period in the EDLP

65 s_1 time period of the off-peak to peak transition time in the EDLP

66 s_2 peak time period in the EDLP

67 s_3 time period of the peak to off-peak transition time in the EDLP

68 t_0 first time interval of the EDLP where the slope of the off-peak/peak tran-
 69 sition starts

70 t_1 first time interval of the EDLP where the main peak stabilises

71 t_2 first time interval of the EDLP where the peak starts to decrease

72 t_3 first time interval of the EDLP where the non-peak behaviour stabilises
 73 after the peak

74 \vec{t} t_0, t_1, t_2 and t_4

75 y year used to compute the EDLP

76 D number of time intervals of the EDLP

77 B set of supermarket building characteristics used to predict the EDLP

78 L_s EDLP of the supermarket s

79 S, S' sets of new and existing supermarkets respectively

80 **2-feat** $\mu(s_0), \mu(s_2), m(s_1), m(s_3)$ and \vec{t}

81 **4-feat** $\mu(s_0), \mu(s_2), m(s_1)$ and $m(s_3)$

82 **8-feat** $\mu(s_0), \mu(s_2), m(s_1), m(s_3)$ and \vec{t}

83 $\mu(s_i)$ mean of the energy values that are in s_i

84 $m(s_i)$ slope of the line that crosses the energy values that are in s_i

85 1. Introduction

86 The aim of reducing greenhouse gas emissions is shared by most coun-
 87 tries [1, 2] with the UK aiming at greenhouse neutrality by 2050 [3]. As energy
 88 use in buildings across the EU accounts for more than 30% of final energy
 89 demand [4], cutting and time-shifting energy demand of all types of buildings
 90 (residential, commercial, services and industrial) are needed to achieve these
 91 targets. Residential buildings have received much attention [5, 6, 7], whilst
 92 commercial and industrial buildings less so due to the lack of open data-sets
 93 and their heterogeneity [8, 9].

94 The total energy demand and the temporal profile are useful performance
 95 indicators for buildings estate management, investment decisions, site acqui-
 96 sitions, and improvement programmes. Knowing the expected demand of a
 97 store establishes a baseline for: 1) planning annual energy budgets for the
 98 portfolio of stores, 2) negotiating energy supply contracts, and 3) detecting
 99 stores with abnormal or anomalous usage. Inevitably there are differences
 100 between stores, with the key being understanding the variability and what
 101 is acceptable usage for any store. Grouping the stores based on common
 102 demand patterns reveals existing distinct behaviours in the store portfo-
 103 lio [10, 11, 12, 13]. This informs which measures might be more effective or
 104 cost-efficient for each group, and identifies stores whose demand patterns do
 105 not match any of the discovered groups (anomalous behaviour).

106 In general, clustering techniques are unsupervised machine learning algo-
 107 rithms that divide data-sets into groups (clusters) without a priori informa-
 108 tion [14, 11]. Both prediction and clustering of energy demand is commonly
 109 performed over electricity daily profiles (EDLP), which are a concise, infor-
 110 mative, and intuitive way to represent, analyse and visualise the electricity
 111 demand of any source [11, 12]. EDLPs are data representations for which
 112 the electricity demand during a day is computed with a temporal granularity,
 113 D . This temporal resolution indicates the number of points (demand values)
 114 that formed the profile, *e.g.* if $D = 24$ each demand value is the hourly
 115 demand. A disadvantage of using full EDLPs is the so-called “curse of di-
 116 mensionality” [15], meaning that some machine learning (ML) algorithms
 117 can have temporal and memory issues when working with high-dimensional
 118 data-sets. By the same token, too few data points for training the algorithms
 119 risk over-fitting the model.

120 We propose representing EDLPs using only a small set of characteristic
 121 features (dimensional reduction) that are automatically extracted from the
 122 profile instead of using the D -dimensional daily profile. We separately in-
 123 vestigate both predicting and clustering the EDLPs using only the extracted
 124 features. New supermarkets profiles are predicted using the historical de-
 125 mand of other stores. Four different ML regression algorithms are imple-
 126 mented to predict EDLPs over a data-set of 213 supermarkets during a six
 127 year period. Experiments in clustering EDLPs are performed using six dif-
 128 ferent algorithms over the supermarket data-set and a data-set of 641 retail
 129 stores, independently. Both data-sets are real data obtained from smart me-
 130 ters. Based on the proposed extracted features and these two ML problems,
 131 the questions that we try to answer are:

- 132 • How accurately can D-dimensional EDLPs be represented using a small
133 set of features?
- 134 • Using only this set of features is it possible to predict future EDLPs
135 of new stores with different ML methods, as accurately as when using
136 the whole EDLP?
- 137 • Using only this set of feature is it possible to cluster the electricity
138 demand as accurately as when using the whole EDLP?
- 139 • Is it possible to extend and generalise this representation over other
140 commercial data-sets that have other temporal resolution?

141 The paper is structured in the following way. We review the literature for
142 ML related to energy analytics in Section 2. In Section 3 we explain the pre-
143 processing of the real-world data-sets and the computational experiments.
144 The results obtained for these experiments and discussion about them are in
145 Section 4. Finally, we draw conclusions and propose future work in Section 5.

146 2. Literature review

147 Predicting electricity demand of buildings (regardless of type) can be di-
148 vided into two basic approaches: model-driven and data-based. The model-
149 driven approach uses sophisticated high-resolution engineering methods based
150 on the thermal, energy, and architectural features of the building to simu-
151 late its energy demand. In data-driven approaches, the energy performance
152 of the building is directly modelled with numerical and statistical methods.
153 There are extensive reviews on methods to predict, model and benchmark
154 energy use in buildings [16, 17, 18, 19, 20, 21]. A possible classification of the
155 data-driven techniques used to predict energy demand [20, 21] is: 1) conven-
156 tional statistical techniques, 2) classification-based models, 3) support vector
157 regression (SVR) model, 4) artificial neural networks (ANN), 5) genetic algo-
158 rithms, 6) grey models, 7) fuzzy model and 8) other models (*e.g.* case-based
159 reasoning). Our study exploits the first four classes of techniques, and we
160 focus our review on the prediction of demand in commercial buildings and
161 supermarkets.

162 Conventional statistical techniques include change-point algorithms and
163 linear regression models such as autoregressive models and ordinary least
164 squares (OLS). Autoregressive models have been used to predict short-term

165 heat load for a single building [22] and, in combination with ANN, used
166 to predict the annual electricity consumption of 787 education facilities in
167 South Korea over a period of seven years [23]. Schrock and Clarige [24] used
168 a change-point algorithm and a year of 15-min electricity readings of one
169 grocery store to predict hourly and daily consumption. Linear regression
170 has been applied to the prediction of 1-h heat load profiles of 116 buildings
171 (health, education, business, and hotels) over three years [25]. The same
172 linear models have also been used on data from 215 UK large supermarkets
173 to estimate the total annual electricity demand [26], and by [27] to estimate
174 annual energy-use intensity for UK 30 supermarkets using building features
175 such as floor area and building age and the number of customers. In the
176 context of climate change adaption [28] exploited temperature and humidity
177 values to predict weekly electricity and gas demand for a single supermarket
178 for the period 2030-2059 using multiple linear regression analysis.

179 Classification-based models include algorithms that were extended to per-
180 form regression. The k-nearest neighbour regression (kNNR) algorithm was
181 used to forecast the next day consumption of 6,000 domestic Irish build-
182 ings in [29], and for the hourly air conditioning load of an office building in
183 China [30]. Random forest (set of decision trees) and ANN (separately) were
184 used to predict hourly HVAC loads of a Spanish hotel [31] over a period of 15
185 months. Similarly, decision trees, ANN and linear regression are compared
186 to predict weekly electricity consumption of 1200 dwellings during the winter
187 and summer of one year [32].

188 ANN has been used to predict the energy demand in 17 studies from 1996-
189 2015 [19]. Short-term electricity demand of a commercial building complex
190 using 15-min resolution data was predicted using ANN in [33]. Daily diurnal
191 cooling load is forecasted for three university buildings with ANN in [34],
192 using data recorded over two years. Both ANN and SVR were compared
193 when predicting hourly cooling load in an office building in China [35] and
194 hourly energy consumption of an office building in Shanghai [36]. Electric-
195 ity consumed by the HVAC and refrigeration systems of one supermarket is
196 predicted using ANN [37]. ANN, Gaussian process regression, linear regres-
197 sion and dynamic mode decomposition are compared in the prediction of 1-h
198 weekday profiles of a commercial building [38]. Lastly, deep learning models
199 (large neural-networks) have been also explored for this problem, however
200 they need large data-sets to estimate the model parameters. For example,
201 Hafeez *et al.* [39] used deep learning for short-term load forecasting over three
202 power-grids with hourly resolution. A deep learning network and a genetic

203 algorithm were combined to predict the 1-h daily profile in an office building
204 over one year [40]. This work applies the clustering of daily weather profile
205 before predicting the demand.

206 Support Vector regression (SVR) models were used by Dong *et al.* [41] to
207 predict monthly energy consumption of four commercial buildings in Singa-
208 pore. Models based on SVR have also been used to predict the energy load
209 (hours to days) of a French residential building [42]. SVR and six other tech-
210 niques was also investigated by [43] to predict next-hour residential building
211 electricity consumption for three houses. Jain *et al.* [44] examine the impact
212 of temporal (*e.g.* daily, hourly, 10 min intervals) and spatial (*e.g.* , whole
213 building, by floor, by unit) granularity to short-term prediction. Experi-
214 ments were performed using SVR over data from a multi-family residential
215 building in the USA. Granell *et al.* [13] compared four techniques, kNNR, or-
216 dinary least of squares linear regression, ANN, and SVR in predicting whole
217 EDLPs of new supermarkets using data from a portfolio of 213 UK super-
218 markets with readings spanning six years.

219 From this range of techniques we can conclude that there is no con-
220 sensus about the superiority of a specific technique. Studies that compare
221 several techniques usually report marginally differences in the prediction re-
222 sults *e.g.* [30, 32, 13, 43], or contradictory results *e.g.* ANN over-performs
223 SVR [36] and vice-versa [35]. These results support our selection of four dif-
224 ferent types of predictors to address our problem. In addition, our prediction
225 work addresses some of the areas that have received less attention. First, re-
226 tail is clearly under-represented in the literature. For example, according to
227 reviews by Chung [16] and Li *et al.* [21] only 22% and 33%, respectively, of
228 investigations were about consumption in commercial buildings, and fewer
229 still in other studies [17, 18]. Particularly notable is the severe lack of work
230 in the literature on predicting energy use by supermarkets. There are differ-
231 ences between patterns of energy demand in commercial and retail premises,
232 but also similarities in niche sectors [8]. Secondly, the prediction of daily pro-
233 files [39, 25, 13, 40] is not common, most of the long-term prediction studies
234 use weekly, monthly, or annual demand. Thirdly, prediction experiments
235 using retail data-sets with a size that can be considered representative (hun-
236 dreds of buildings) are also infrequent. Finally, predicting the future demand
237 of new buildings for a long period of time (more than three or four years) is
238 a highly unusual approach; most studies predict the future demand for the
239 study building and they usually do not use several years of continuous data.

240 Reviews of clustering methods applied over electrical data can be found

241 in [11, 12, 10]. Most studies have used residential data-sets, but some work
 242 clustering electricity profiles of commercial and industrial customers has been
 243 completed. For example, 292 Greek industrial and service customers are
 244 clustered using a two-stage ML algorithm [45]. Wavelet decomposition was
 245 used [46] to select significant features to describe the hourly load profiles
 246 of 9,092 Danish industrial and commercial loads for two-week data. Later,
 247 they applied clustering using the k-means algorithms over these features.
 248 In [47] they investigate several clustering techniques such as k-means and
 249 hierarchical algorithms to cluster 234 non-residential customers, and a data
 250 set of 1,877 UK business from the entertainment sector was used to perform
 251 clustering with a Dirichlet process mixture model [48].

252 A recent review of dimensional reduction techniques appears in [10]. Di-
 253 mensional reduction has been attempted for electricity demand modelling
 254 and clustering [46], and for symbolic aggregate approximation with hier-
 255 archical clustering [49]. Representing the data with principal component
 256 analysis, the curvilinear component analysis, and the Sammon map are in-
 257 vestigated by [47]. The effect of the time resolution when clustering domestic
 258 EDLPs [50] was investigated by averaging over regular intervals instead of
 259 extracting key features based on the specific shape of the retail EDLP as we
 260 do here. Residential demand profiles have been characterised and clustered
 261 with a set of five points that match the peaks [51].

262 3. Methods

263 First we describe the data-sets used to perform the experiments. Sec-
 264 ondly, the features to represent the EDLP and methods to extract them are
 265 explained. Thirdly, we describe the prediction algorithms, evaluators, and
 266 experiments. Finally, clustering algorithms and evaluators are defined.

267 3.1. The data-sets

268 Two data-sets are used to perform the experiments. The first comprises
 269 1-h resolution electricity meter readings (kWh) from 213 UK supermarkets
 270 of the same chain for the period 2012–17. The detail of the meta-data fea-
 271 tures available of each supermarket are described elsewhere [13], but are
 272 summarised as: floor area subdivided into eight categories (*e.g.* chilled, pro-
 273 duce, storage), geographical location, daily average external temperature,
 274 and electricity consumption. There are 129 supermarkets that use electricity

275 and gas (SEG) and 84 supermarkets that use only electricity (SE). The sec-
 276 ond data-set comprises 663 UK retail stores (single company) with electricity
 277 meter readings at 0.5-h resolution acquired between April 2013 and October
 278 2014. In this case, the only meta-data fields are the address and outlet type
 279 category that summarise the location of the store (*e.g.* arterial route, high
 280 street retail park, shopping centre).

281 For both data-sets an analytic filtering pre-process removes anomalous
 282 readings with zero or negative values, accounting for less than 0.8% of the
 283 data. In addition, stores with less than the equivalent of half a month of
 284 data (360 and 720 readings for the supermarket and retail store data-sets
 285 respectively) are removed: For the retail store data-set this was 22 shops
 286 leaving 641 stores for analysis, whilst for the supermarkets it varied from
 287 year-to-year [13].

288 3.2. Features extraction to represent the EDLP of supermarkets

289 Like most retailers, the supermarkets have a fixed daily schedule: they
 290 usually open in the morning to close later in the evening [52]. Based on
 291 these schedules, the electricity consumption patterns are quite similar to each
 292 other with a typical inverted-U shape. Figure 1 shows the daily profiles, for
 293 different seasons, of four different supermarkets and retail stores from our
 294 data-sets. These eight EDLPs show similar patterns characterising the peak
 295 and off-peak periods, however, they exhibit variability during these periods.
 296 Energy demand by supermarkets are greater than that of retail stores.

297 Based on these behaviours we can define four time periods in which im-
 298 portant changes occur (Figure 2):

299 t_0 indicates the first time interval of the EDLP where the slope of the off-
 300 peak/peak transition starts.

301 t_1 is the first time interval of the EDLP the main peak stabilises.

302 t_2 is the first time interval of the EDLP the peak starts to decrease.

303 t_3 is the first time interval of the EDLP where the non-peak behaviour sta-
 304 bilises after the peak.

305 These periods follow the conditions that $t_i \in [0, D - 1], 0 \leq i \leq 3$ and
 306 $t_i < t_{i+1}, 0 \leq i \leq 2$. In the example given in Figure 2 their value are: $t_0 = 6$,
 307 $t_1 = 9$, $t_2 = 15$ and $t_3 = 21$, corresponding to 6.00am, 9.00am, 3.00pm and

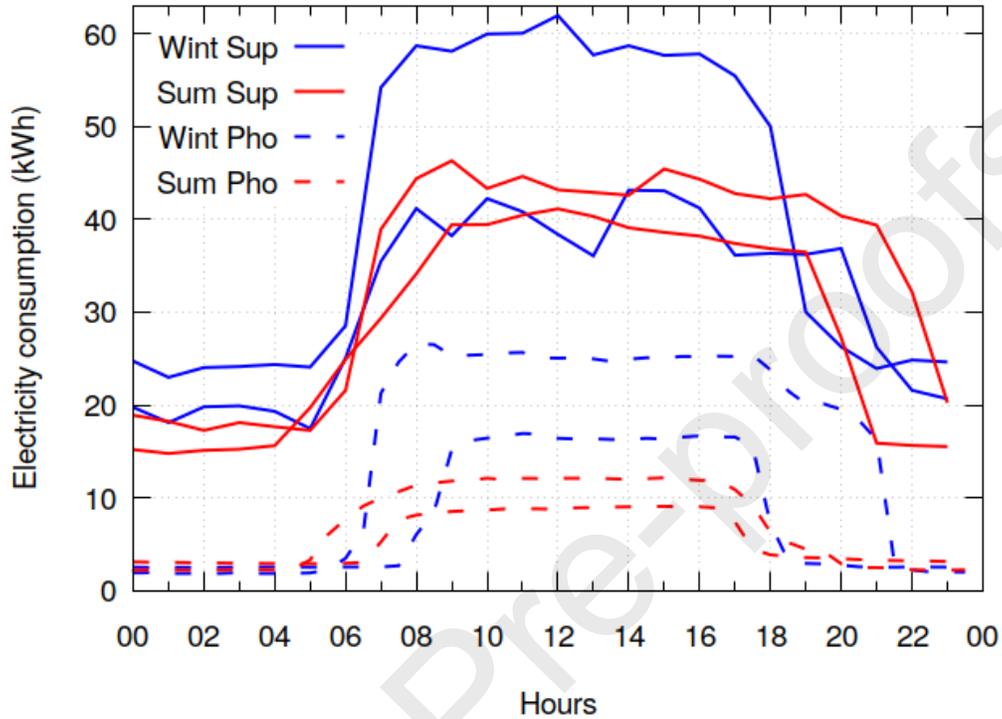


Figure 1: Example Winter and Summer daily profiles of four different supermarkets (Sup) and retail stores (Pho).

308 9.00pm, respectively. By defining the vector grouping the four time features
 309 as $\vec{t} = \{t_0, t_1, t_2, t_3\}$, we can divide the EDLP into four intervals using:

310 **off-peak** time period in which the supermarket is closed and the demand
 311 is a stable baseload of refrigeration, as HVAC and lighting should
 312 be switched-off or to minimum power. Formally, it is $s_0 = [0, t_0 -$
 313 $1] \cup [t_3, D - 1]$ *e.g.* horizontal green lines in Figure 2.

314 **off-peak to peak transition** short period occurring a little before the store
 315 is opened to customers when the HVAC, lighting, and other services
 316 ramp to their peak values. Formally, it is $s_1 = [t_0 - 1, t_1]$ *e.g.* horizontal
 317 yellow line in Figure 2.

318 **peak** period in which the demand is constantly high as the supermarket is
 319 open. The appliance power consumption is usually stable, but short-

320 term variability may occur (see EDLPs of Figure 1). Formally, it is
 321 $s_2 = [t_1, t_2 - 1]$ *e.g.* horizontal pink line in Figure 2.

322 **peak to off-peak transition** short period following the closure of the store
 323 to customers, but staff may still be present. Modern appliances should
 324 not have a very long temporal lag for reducing their demand when they
 325 are switched-off. Formally, it is $s_3 = [t_2 - 1, t_3]$ *e.g.* horizontal grey
 326 line in Figure 2.

327 Given any interval of time $s = [t, t']$ with $t' > t$, we define two generic
 328 operators: 1) $\mu(s)$ as the mean of the energy values from time t to t' ,
 329 *i.e.* $\mu(s) = \sum_{i=t}^{t'} e_i / (t' - t + 1)$ and 2) $m(s)$ is the slope of the line that
 330 crosses the points (t, e_t) and $(t', e_{t'})$, *i.e.* $m(s) = (e_{t'} - e_t) / (t' - t)$.

331 We can describe the profile using eight features: the four time periods
 332 of the events (\vec{t}), consumption of the off-peak and peak periods ($\mu(s_0)$ and
 333 $\mu(s_2)$), and the slopes of the transitions ($m(s_1)$ and $m(s_3)$). The demand
 334 values of $\mu(s_0)$ and $\mu(s_2)$ are the average during all the values of the off-peak
 335 and peak respectively, and they are a linear approximation of the demand
 336 during these time intervals. Values of $m(s_1)$ is the rate of demand increasing
 337 by hour when moving from off-peak to peak period (this value is always
 338 positive as demand increases during this period.). The value of $m(s_3)$ is
 339 always negative as the demand decreases during the peak/off-peak transition
 340 interval. Given these eight features, the estimated profile $\vec{e}' = \{e'_0, \dots, e'_{D-1}\}$
 341 can be reconstructed using Euclidean geometry:

- 342 • Off-peak values are equal to $\mu(s_0)$: $e'_i = \mu(s_0), 0 \leq i < t_0$ and $t_3 \leq i <$
 343 D
- 344 • Values of the off-peak/peak transition are computed with the linear
 345 equation $y = x * m(s_1) + b$ where independent term b is computed
 346 by substituting the equation with the data point $(t_0 - 1, \mu(s_0))$: $e'_i =$
 347 $i * m(s_1) + b, t_0 \leq i < t_1$.
- 348 • Peak values are equal to $\mu(s_2)$: $e'_i = \mu(s_2), t_1 \leq i < t_2$.
- 349 • Values of the peak/off-peak transition are calculated with the linear
 350 equation $y = x * m(s_3) + b'$ where term b' is computed by substituting
 351 equation with the data point $(t_2 - 1, \mu(s_2))$: $e'_i = i * m(s_3) + b', t_2 \leq i < t_3$.

352 As an example, Figure 2 shows a reconstructed profile (red lines) obtained
 353 using the eight selected features and a real profile (black line) of a supermar-
 354 ket EDLP. The off-peak demand is estimated well, likewise the central part
 355 of the peak demand. However, the beginning of the peak demand is underes-
 356 timated and the end is overestimated. The discrepancy (error) between the
 357 reconstructed profile ($\vec{e}' = \{e'_0, \dots, e'_{D-1}\}$) and the real values of the profile
 358 ($\vec{e} = \{e_0, \dots, e_{D-1}\}$) is quantified using evaluators:

359 **Euclidean Distance (ED)** in which discrepancies between the EDLPs ab-
 360 solute values are accumulated (in kWh),

$$\sqrt{\sum_{i=0}^{D-1} (e_i - e'_i)^2} \quad (1)$$

361 **Normalised Percentage (NP)** difference with respect to the original EDLP
 362 (NP) computes the relative distance considering the proportion of the
 363 error with respect to the total consumption of the original profile,

$$\frac{100 * \sum_{i=0}^{D-1} |e_i - e'_i|}{\sum_{i=0}^{D-1} e_i} \quad (2)$$

364 The ED and the NP between the modelled and real EDLPs of Figure 2 are
 365 15.8 kWh and 3.9% respectively. The evaluators \overline{ED} and \overline{NP} are extended
 366 over the whole data-set using the average ED and NP respectively for all
 367 stores.

368 As the whole feature set can be obtained directly with the time period
 369 vector \vec{t} , they can be automatically computed searching using the objective
 370 function to minimise the error:

$$\hat{\vec{t}} = \arg \min_{\vec{t}} (\text{Ev}(\vec{e}, \vec{e}'_{\vec{t}})) \quad (3)$$

371 where $\vec{e}'_{\vec{t}}$ is the reconstructed profile using \vec{t} , and Ev is an evaluator computed
 372 over the two EDLPs. For evaluator Ev, we use the ED. A brute-force search
 373 method in which all possible values of \vec{t} are explored to find the optimal
 374 solution $\hat{\vec{t}}$ as it is restricted search. For the example (Figure 2) the set of
 375 features obtained using this objective-function method are $\vec{t} = (6, 9, 15, 21)$,
 376 $\mu(s_0) = 32.0$ kWh, $\mu(s_2) = 100.0$ kWh, $m(s_1) = 16.2$ kWh/h and $m(s_3) =$
 377 -9.2 kWh/h. The utility of this approach needs to be demonstrated for
 378 problems such as prediction and clustering.

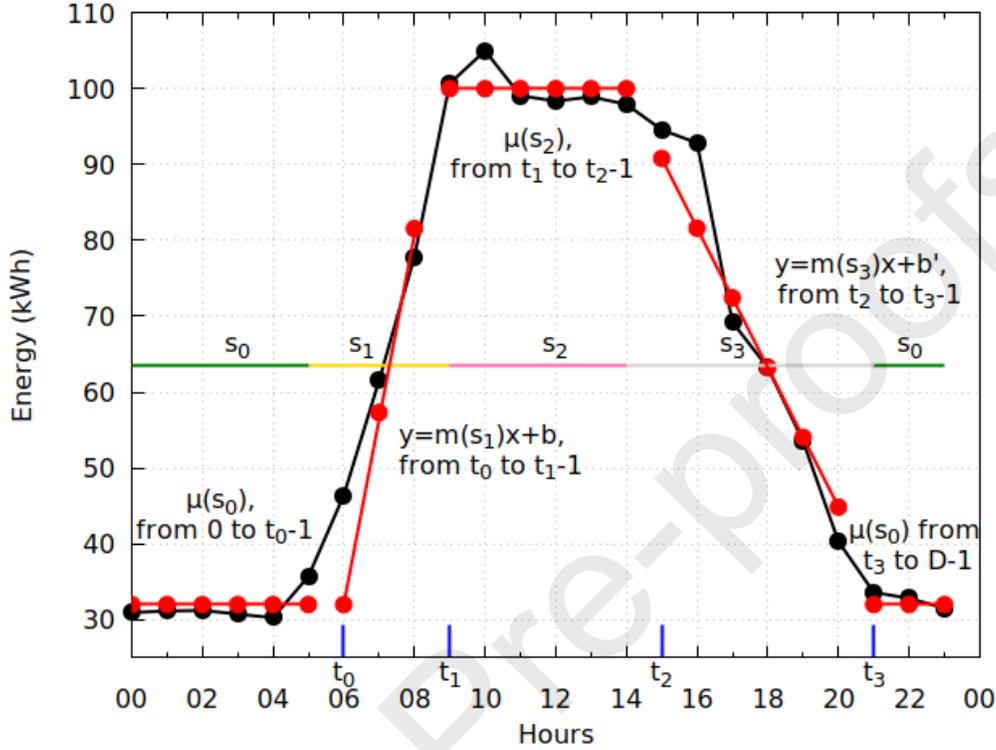


Figure 2: Modelled profile based on the eight proposed features (red line) and real profile (black line).

379 3.3. Computational prediction experiments

380 Experiments are performed using only the extracted features to predict
 381 electricity demand of new supermarkets. The EDLP of a new supermarket
 382 $L_s = e_0, \dots, e_{D-1}$ for a year y is predicted based on historical profiles of
 383 existing supermarkets S' and the supermarket building characteristics B . L_s
 384 is the EDLP of the new supermarket s , e_i is the electricity consumed (kWh)
 385 between the $(i-1)$ -th and i -th time interval, D is the number of intervals, S
 386 and S' are the set of new and existing historical supermarkets, respectively
 387 ($S \cap S' = \emptyset$). The set of store characteristics B is the set of available informa-
 388 tion about the supermarket building such as the floor area divided by usage
 389 and the supermarket geographical location. Therefore, we train a regression
 390 ML algorithm with all the supermarkets of S' where the independent vari-
 391 ables (input) are the store characteristics B and the dependent variable to
 392 predict (output) is e_i , $0 \leq i < D$.

393 As we do not know which store characteristics B use nor how many stores
 394 k to select to train the ML model, the best combination of (k, B) is searched
 395 using Equation 4.

$$(\hat{k}, \hat{B}) = \arg \min_{k, B} \sum_{s \in S} \text{Ev}(L_s, L_s(k, B)) \quad (4)$$

396 where S is the set of new supermarkets, L_s is the real EDLP of supermarket
 397 s , $L_s(k, B)$ is the predicted energy profile when using parameters (k, B) and
 398 $\text{Ev}(L_s, L'_s(k, B))$ is the evaluator that measures the error between the pre-
 399 dicted and real profile. As Ev we use the average Euclidean distance over all
 400 the real and predicted stores:

$$\overline{ED} = \frac{\sum_{s \in S} ED_s}{|S|} \quad (5)$$

401 where ED_s is the ED computed over the real and predicted EDLPs of the
 402 supermarket s .

403 Four different ML algorithms are investigated: kNNR [14], ANN [14], SVR [53],
 404 and OLS [54].

405 We only use the extracted features to represent the EDLP *i.e.* these fea-
 406 tures are predicted using as input the store characteristics B instead of pre-
 407 dicting the whole profile. The diagram of Figure 3 illustrates the steps of the
 408 experimental set-up:

- 409 1. The eight features of each supermarket (\vec{t} , $\mu(s_0)$, $\mu(s_2)$, $m(s_1)$ and
 410 $m(s_3)$) are computed.
- 411 2. These features are predicted independently for each supermarket s'
 412 using the regression model using as input the store features (B'_s). That
 413 is, for each supermarket s' , the eight features of the EDLP of year y
 414 are predicted with the regression algorithm. This ML model is trained
 415 with the features extracted of the EDLP computed in previous years
 416 $y - t$ of the stores of the set $S - \{s'\}$.
- 417 3. The profile of the predicted store is reconstructed with the eight pre-
 418 dicted features of the store (\vec{t}' , $\mu'(s_0)$, $\mu'(s_2)$, $m'(s_1)$ and $m'(s_3)$). The
 419 evaluators are computed between this reconstructed profile and the
 420 original profile of the test supermarket (s').
- 421 4. Parameter search (k, B) is performed and final error is computed over
 422 the best parameter combination (\hat{k}, \hat{B}) that minimizes Equation 4.

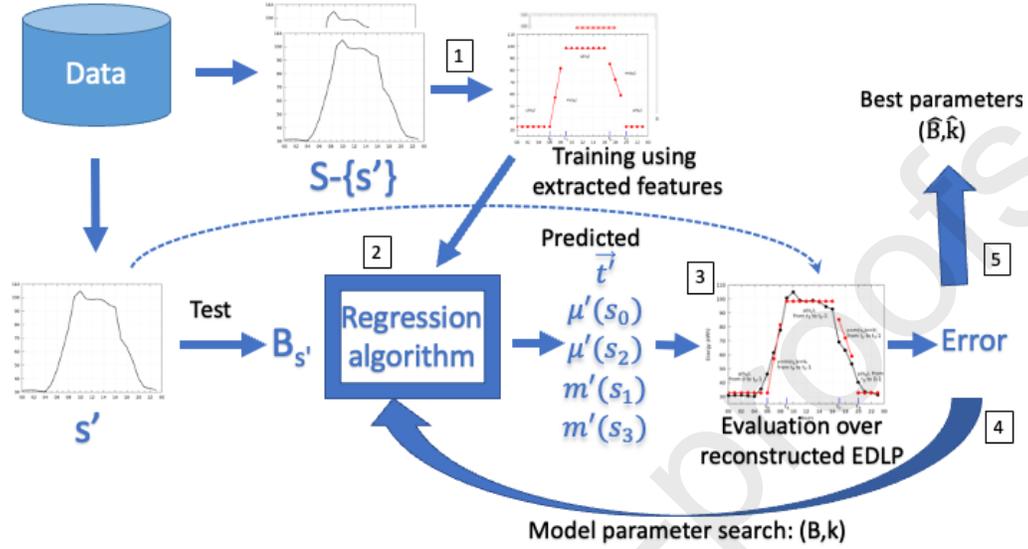


Figure 3: Logical flow of the prediction experiments using the features to represent the profiles.

423 The two essential points of this experimental set-up are 1) the ML algo-
 424 rithm predicts the summarised features of the profile, and 2) the evaluation
 425 is performed comparing the reconstructed profile using the predicted features
 426 with the real profile to predict (not with the reconstructed profile over the
 427 real features). Due to the second point it is feasible to compare the results
 428 obtained with these experiments with the results obtained when predicting
 429 the whole profile. As the values of \vec{t} are integer numbers, the closest integer
 430 is selected to the value returned by the regression model.

431 Fewer than 30 supermarkets are opened each year and we assume that
 432 each is opened within year y . The historical EDLPs of the other $|S| - 1$
 433 supermarkets are used to predict the EDLPs of the new ones, improving the
 434 robustness of the experiments. The leaving-one-out technique [14] uses all
 435 the data points — except the one being estimated — as predictors (repeated
 436 $|S|$ times) to compute the EDLP of the new one for year y .

437 Experiments are carried out separately over EDLPs of the supermarkets
 438 computed for 2013–2017, seasons (Winter, Summer and Spring/Autumn),
 439 and SE/SEG sets [13]. We employ the brute-force approach (Equation 4) to
 440 search all parameter combinations (\hat{k}, \hat{B}) . The maximum number of combi-
 441 nations, for each season is $(2^{|F|} - 1) * (|S| - 1) = (2^{11} - 1) * (129 - 1) = 262,016$.

442 For ANN and SVR (temporally more complex) we used stepwise regres-
 443 sion [14] with the whole feature set B (using all the supermarkets, $k = |S|$).
 444 For the ANN, we use a logistic function as an activation function, over a two
 445 internal layer net, *i.e.* the configuration of the network is $|B|$ -4-2-1, where
 446 $|B|$ is the number of features. The function neuralnet of the R language [55]
 447 is used with the default parameters, *i.e.* the resilient backpropagation algo-
 448 rithm with 10^5 maximum steps for the net training. For SVR, we used a
 449 radial basis kernel function to model the non-linearly. The function svm of
 450 the R language [56] is used with default parameters. The parameters of the
 451 ML methods are the same that were used in [13] to enable comparison with
 452 previous work. Both R scripts were invoked for each one of the computing
 453 experiment from the generic C++ code.

454 3.4. Clustering experiments

455 Clustering experiments group all the available EDLPs computed dur-
 456 ing a specific year for each data-set independently. The result depend on
 457 both the algorithm and the way the data is represented. Our aim is to
 458 compare clustering results—not algorithm performance—with the two data
 459 representations. Thus we selected two types of clustering algorithm: parti-
 460 tioning and agglomerative hierarchical. The partitioning algorithm we chose
 461 was k-means [47, 11, 45, 57, 58]. For the agglomerative hierarchical algo-
 462 rithm [47, 11, 45, 57] there is more choice depending on the criterion used to
 463 compute the distance to merge the clusters: Single link algorithm, Complete
 464 link algorithm, Unweighted pair group method average algorithm (UPGMA),
 465 Unweighted pair group method centroid algorithm (UPGMC), Weighted pair
 466 group method centroid algorithm (WPGMC) and Ward or minimum variance
 467 algorithm (WARD).

468 We selected six evaluators [59] to asses the clustering results: the cluster-
 469 ing dispersion indicator (CDI), Davies-Bouldin index (DBI), modified Dunn
 470 Index (MDI), mean index adequacy (MIA), scatter index (SI), and variance
 471 ratio criterion (VRC). These evaluators are based on the similarity of the
 472 data elements within each cluster, and the difference among elements of the
 473 other clusters.

474 The clustering is performed using directly three sets of features:

475 **8 features (8-feat):** $\mu(s_0)$, $\mu(s_2)$, $m(s_1)$, $m(s_3)$ and \vec{t} .

476 **4 features (4-feat):** $\mu(s_0)$, $\mu(s_2)$, $m(s_1)$ and $m(s_3)$.

477 **2 features (2-feat):** $m(s_1)$ and $m(s_3)$.

478 However, we decided to evaluate directly over the whole profiles. The reason
 479 for this is that the output of clustering is the grouping in which all the data-
 480 points (in our case EDLP) can be distributed based on the ML algorithm.
 481 As all the evaluators use the intra-point distance, we consider that the fairest
 482 way to compare the quality of the obtained grouping is to compare over the
 483 same set of points. Clustering results using the eight features are compared
 484 with respect to the clustering obtained using the whole EDLP. For the k-
 485 means algorithm 100 repetitions with different random initialisation were
 486 performed and the evaluations are averaged. The number of clusters (input
 487 parameter of the algorithm) is 2–10 exploring all the values. All the software
 488 was coded in C++.

489 **4. Results and Discussion**

490 We have performed a large number of computational experiments. For
 491 clarity, we discuss separately the results obtained for: 1) representing the
 492 EDLPs using the eight features, 2) the prediction experiments and 3) the
 493 clustering experiments. Prediction experiments were not performed using
 494 the retail stores data-set as there was only one year of data.

495 *4.1. Representing supermarket EDLPs with the selected features*

496 An example of the features for the Winter 2017 profile of a SEG super-
 497 market is in Section 3.2. Histograms of Figure 4 show the range of values for
 498 the features t_0 , t_2 , $m(s_1)$ and $\mu(s_2)$ extracted from the Winter 2017 profiles
 499 of all the 129 SEG supermarkets.

500 For the periods t_0 (Figure 4a) and t_1 , there are only four different hours
 501 in which they occur, and one of the hours is much frequent than the oth-
 502 ers: 6am (70.5% of supermarkets) and 8am (50.4%) for t_0 and t_1 respectively.
 503 The period t_3 also has one value more frequent than the others (9pm, 55.0%),
 504 however there are eight different values for the t_2 (Figure 4b). The histograms
 505 exhibit little variability of values and the distribution is Gaussian. However,
 506 the most important insight is the variability in which the peak and off-peak
 507 can begin and end. This shows that using a fixed time for these moments is
 508 an over-simplification that does not properly represent the real pattern of the
 509 demand. In addition, the range of values for these time slots is restricted,
 510 indicating common patterns for the supermarkets. Intervals for the mean

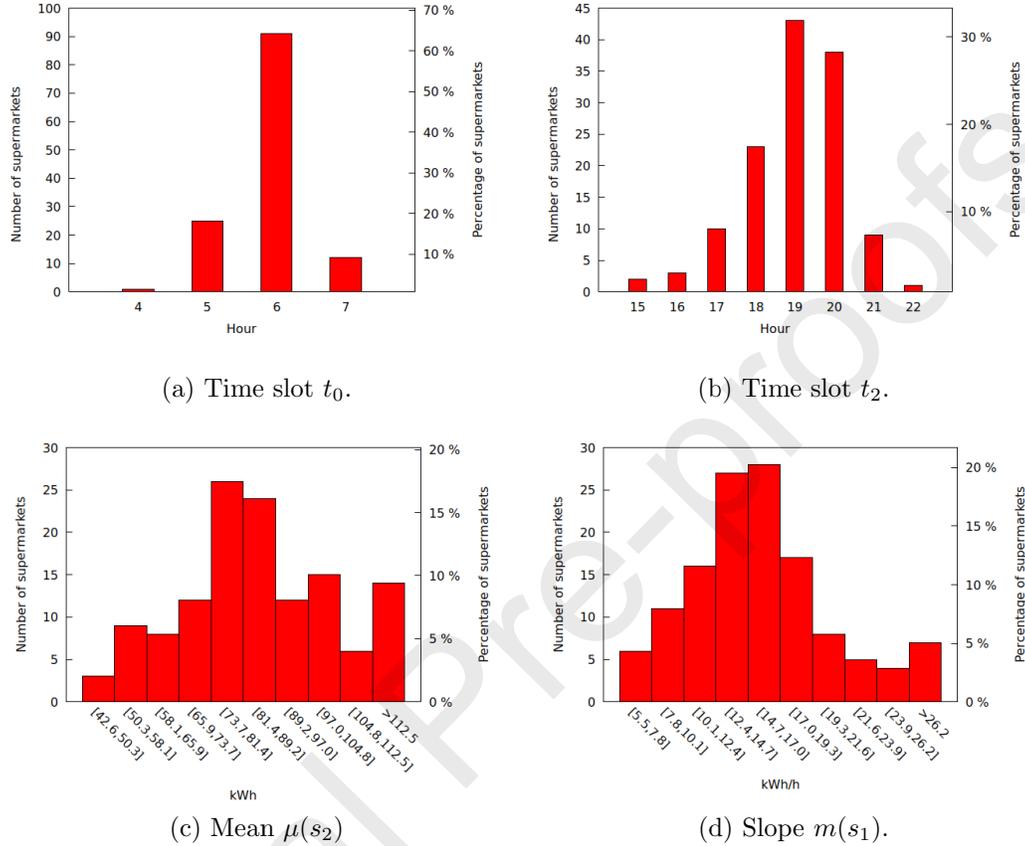


Figure 4: Histograms with values for t_0 , t_2 , $m(s_1)$ and $\mu(s_2)$ features computed over the SEG supermarkets (Winter 2017 profiles).

511 values (kWh) in Figure 4c and slopes values (kWh/h) in Figure 4d need to
 512 be used as they are continuous variables. Nine different intervals are created
 513 for the histograms and an additional bucket with the extreme values. Both
 514 average demand values for peak and off-peak periods show an important vari-
 515 ability in their respective values. One reason for this large range of demand
 516 values is the large variability of the floor area. These two histograms are not
 517 normally distributed.

518 The same analysis can be computed for any season, year and store type.
 519 Results would be similar as the variability of the shape of the profiles is not
 520 very high. This is demonstrated when computing the errors. The \overline{NP} evalu-
 521 ator between the real the reconstructed profiles (all years, seasons and set of

522 stores) are computed. For the SE stores, the best \overline{NP} results are 5.8%, 4.5%
 523 and 5.3% for 2017 Winter, Summer and Spring/Autumn ELDPs respectively.
 524 For the SEG stores, the best \overline{NP} results are 4.3%, 4.1% and 4.1% for 2017
 525 Winter, Summer and Spring/Autumn ELDPs respectively. We note that
 526 the error increases when the profiles are computed over older years. The
 527 worst \overline{NP} scores is 7.2% computed over stores just with electricity over the
 528 Winter 2014 profiles. Comparing seasons, errors over Winter profiles are al-
 529 ways slightly greater than for Spr/Aut profiles and errors over these ones are
 530 greater than for Summer profiles. The error for stores that consume electric-
 531 ity and gas is lower than stores than consume only electricity. This indicates
 532 that the heating system increases the complexity of the profile making the
 533 approximation of it using the proposed features more difficult. In analysing
 534 the shape of the profiles, we see that the demand fluctuations during the
 535 main peak are more common in Winter than Summer profiles, *e.g.* the 10am
 536 peak in Figure 2, or the afternoon in Figure 1. These fluctuations increase
 537 the error when modelling the consumption by averaging the demand over
 538 long periods, as we do with the reconstructed profile.

539 When representing the retail stores using the proposed feature and all
 540 the seasonal data $\overline{ED} = 1.0kWh$ and $\overline{NP}=3.8\%$. Errors for this data-set
 541 are lower than errors obtained with the supermarkets as they have lower
 542 demand and a more regular U-inverted shape. Figure 5 displays the real
 543 and re-computed EDLPs for the case with the lowest NP (0.5%), median NP
 544 (3.5%) and worst NP (11.7%). The reconstructed EDLP in Figure 5a and
 545 Figure 5b match quite well the respective real EDLP. In the case of Figure 5c,
 546 the error is greater as there is an additional peak in the peak period and a
 547 valley in the off-peak period. Our model does not represent properly such
 548 events, but this type of event is unusual. Similar scores can be seen when
 549 using the Summer, Winter and Spring/Autumn profile.

550 4.2. Prediction experiments

551 Prediction experiments are independently performed for all supermarket
 552 EDLPs computed during each year (2013-2017), season (Winter, Summer and
 553 Spring/Atumns) and store type (SE and SEG) giving a total of $5*3*2=30$
 554 different sets. An example of prediction for a particular supermarket (the
 555 example in Figure 2) is shown in Figure 6. The black profile is the real
 556 demand, the red and green profiles are predicted using the feature represen-
 557 tation. These were the best predictions (considering the parameter search
 558 that minimises \overline{ED} over all the set of stores) and they were obtained using

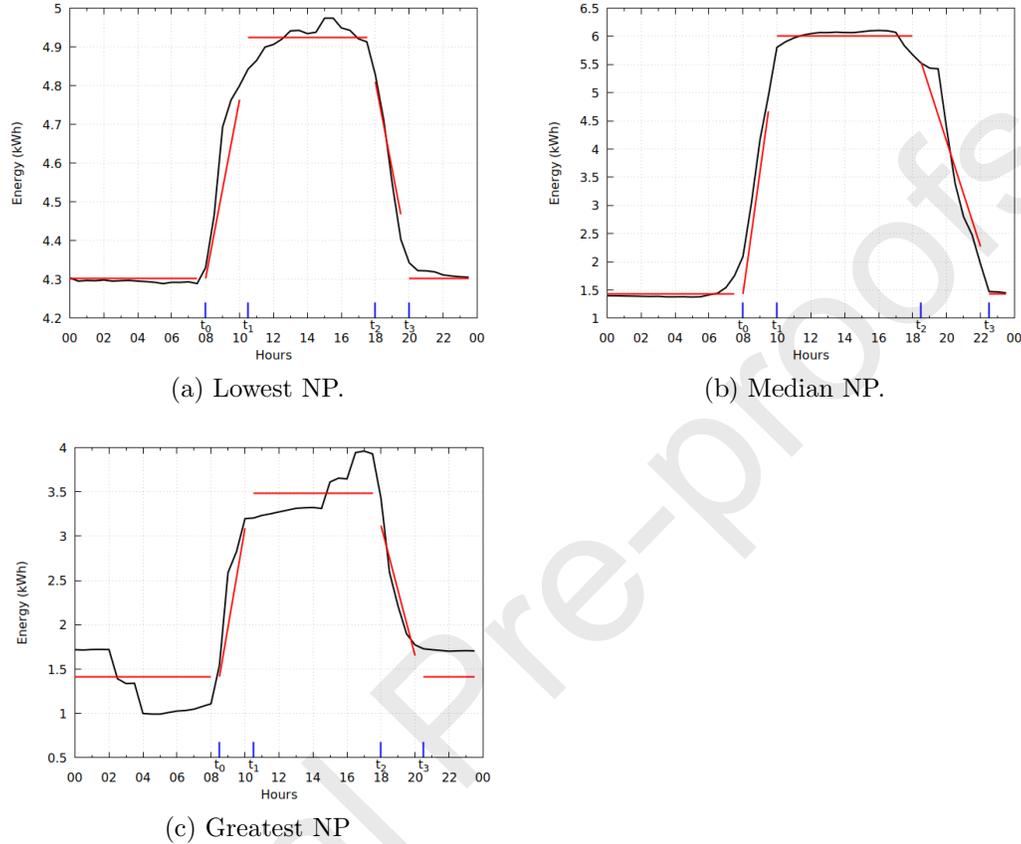


Figure 5: Real and reconstructed EDLP using the features with the lowest, median and worst NP scores for the retail store data-set.

559 OLS with features= $\{\text{GM area, Cafeteria area, Sales area, Office area, Chilled}$
 560 $\text{area}\}$ and $k=98$ for the whole profile representation and features= $\{\text{GM area,}$
 561 $\text{Cafeteria area, Sales area, Storage area, Chilled area, Location}\}$ and $k=75$
 562 for the key feature representation. The values for the evaluators are $\text{ED}=64.0$
 563 kWh and $\text{NP}=16.6\%$ when predicting the features and $\text{ED}=59.5 \text{ kWh}$ and
 564 15.6% when predicting the whole profile. In this case, using the features
 565 implies a relative increase of the error of 7.5% and 6.4% with respect to the
 566 whole profile prediction for ED and NP evaluators, respectively.

567 Table 1 shows the results for the \overline{NP} evaluator obtained when averaging
 568 the evaluator over all the supermarkets in the set.

569 The lowest error for the \overline{NP} evaluator is 12.5% (Summer 2017) using

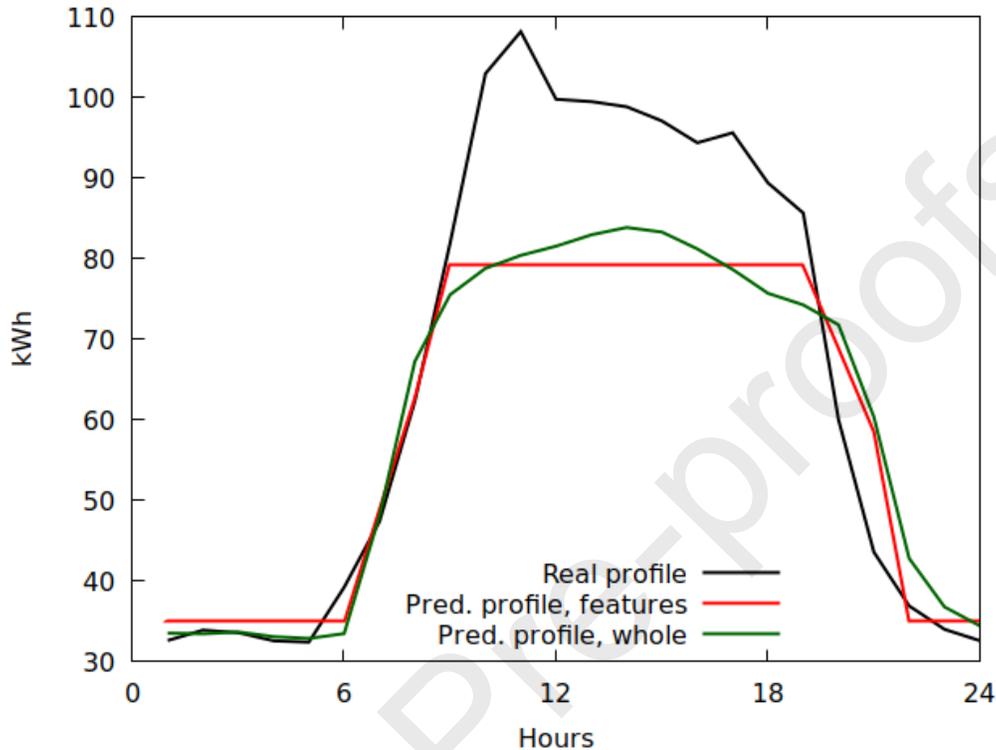


Figure 6: Real and predicted profiles using both the feature representation and the whole profile for one supermarket.

570 the OLS regression method for supermarkets using electricity and gas. This
 571 result is in line with those for the whole profile [13] and can be summarised
 572 thus:

- 573 • Errors computed over cold seasons are greater than errors obtained
 574 during warm seasons *i.e.* Summer profiles are better predicted than
 575 Spring/Autumn profiles, which are better than Winter profiles. The
 576 most likely cause is the uncertainty and variability of the heating system
 577 consumption:
- 578 • Errors obtained during most recent years are usually smaller than for
 579 old data. We suggest that stores tend to become more homogeneous
 580 as older appliances are routinely replaced.
- 581 • There are only small differences when comparing algorithms. However,

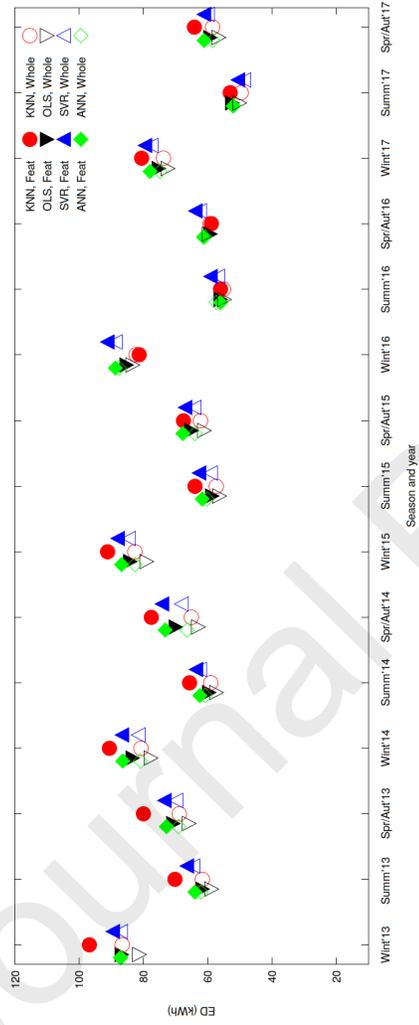
582 the OLS usually outperforms the other three regression methods which
583 is due to the modest size of the data-sets.

- 584 • Stores with electricity and gas are better predicted than stores using
585 electricity only. This too relates to the level of complexity added by
586 the need to also predict the heating demand.

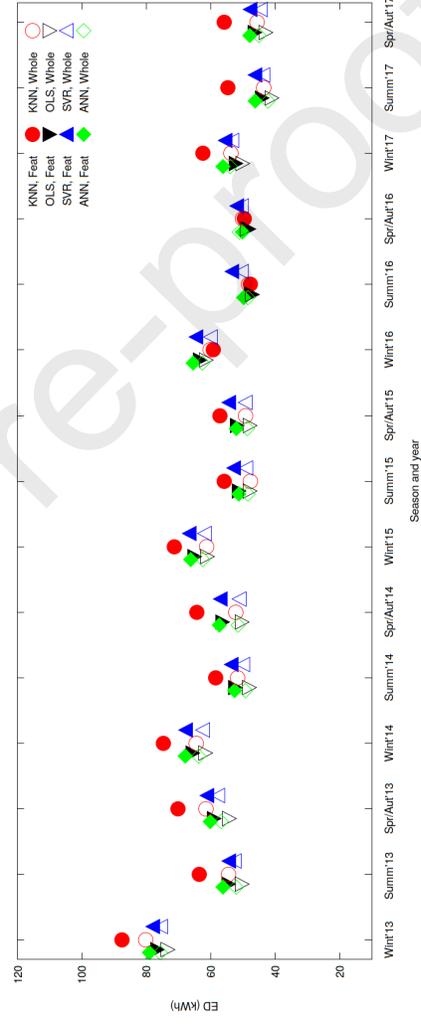
587 Comparing the results obtained using the feature set and those using
588 whole profile representations shows the feasibility of exploiting reduced di-
589 mensionality to predict EDLPs. Figure 7 shows the \overline{ED} values using both
590 representations. The scores when using the full dimensional set (the whole
591 profile) to predict the ELDP are better than using the reduced feature set.
592 However, in many cases the difference is insignificant, especially for the most
593 recent years. Using the \overline{ED} evaluator the absolute difference is an average
594 of 4.0 kWh (6.0%) and 4.4 kWh (8.3%) for SE and SEG, respectively, when
595 comparing the two methods. For both SE and SEG, \overline{NP} using the feature
596 set is 0.9 points worse than using the whole profile. The relative differences
597 for this evaluator are 4.6% and 5.9% for SE and SEG respectively.

TypSt	Year	Season	KNN	OLS	SVR	ANN
Stores just with elec. (SE)	2013	Wint	23.5	22.0	21.2	22.5
		Sum	20.8	18.9	19.4	19.6
		Spr/Aut	22.1	19.3	19.4	20.3
	2014	Wint	23.2	21.9	22.6	23.2
		Sum	20.6	19.2	20.2	20.5
		Spr/Aut	24.9	21.4	22.9	22.4
	2015	Wint	25.1	22.7	23.9	23.3
		Sum	23.0	20.2	20.9	21.5
		Spr/Aut	21.8	20.6	20.9	21.4
	2016	Wint	25.2	26.3	27.9	27.4
		Sum	19.7	18.6	18.8	19.6
		Spr/Aut	19.0	19.0	19.6	20.3
	2017	Wint	22.9	21.9	22.8	23.0
		Sum	17.7	18.1	17.6	19.2
		Spr/Aut	21.2	19.6	19.9	20.2
Stores with elec. and gas (SEG)	2013	Wint	21.5	18.5	18.9	19.3
		Sum	16.3	13.9	13.9	14.3
		Spr/Aut	17.9	15.2	15.8	15.5
	2014	Wint	19.9	17.1	17.9	18.6
		Sum	16.3	14.9	14.9	14.9
		Spr/Aut	17.3	15.6	15.9	15.8
	2015	Wint	18.7	17.4	17.9	17.9
		Sum	16.1	15.0	15.5	15.1
		Spr/Aut	16.2	14.7	15.6	15.3
	2016	Wint	17.2	17.7	18.1	18.6
		Sum	13.6	13.1	14.9	13.7
		Spr/Aut	14.3	13.5	14.4	14.1
	2017	Wint	17.5	14.6	15.6	16.2
		Sum	15.3	12.5	13.1	13.2
		Spr/Aut	16.0	13.1	13.7	13.9

Table 1: Prediction results using the \overline{NP} (%) evaluator for the profile represented with the key features. Results are separated by algorithms, seasons, years and store type.



(a) Stores just with electricity (SE)



(b) Stores with electricity and gas (SEG)

Figure 7: Prediction results evaluated using the \overline{ED} (kWh) for the profile represented with the reduced feature set (solid colour) and the whole profile. Points are offset in the x-axis for clarity.

598 To understand the reasons for the greater error using reduced dimen-
 599 sionality it is necessary to re-think the sequence of processes performed in
 600 this prediction experiments (Figure 3). In this sequence, both modelling
 601 and prediction errors can occur throughout in the process chain. First, the
 602 profile to be predicted is modelled using the features with non-trivial er-
 603 ror (see Section 4.1). Secondly, like any prediction process the features of
 604 the EDLP are not perfectly estimated using the regression model. Thirdly,
 605 when reconstructing the profile using these predicted features we are again
 606 approximating the whole profile adding new error.

607 As the evaluation is performed against the (full dimensional) real profile
 608 it seems logical to have greater error than predicting directly whole profiles.
 609 On the other hand, we have shown that the features are able to explain and
 610 capture the main patterns of the load profile with fewer parameters to predict
 611 than using the whole profile. Interestingly, as the difference in the results
 612 are small, the positive factors compensate the negative ones indicating the
 613 feasibility of using reduced dimensionality.

614 4.3. Clustering experiments

615 Clustering experiments are performed independently for all supermarket
 616 EDLPs computed during each year, season, and store type (SE, SEG and
 617 both together); a total of $5*3*3=45$ experiments. Figure 8 shows the results
 618 obtained when clustering the EDLPs only represented with $\mu(s_0)$ and $\mu(s_2)$
 619 (2-feat) when using readings during Winter 2017 of SEG supermarkets and
 620 the k-means algorithm ($k=4$). The clusters show a clear separation (Fig-
 621 ure 8a), especially in the $\mu(s_2)$ feature because the value of $\mu(s_2)$ is greater
 622 than $\mu(s_0)$, giving more weight when computing distances among clusters.
 623 The real EDLPs of each cluster are used to compute the evaluators. The
 624 profile of each cluster centroid (Figure 8b) are distinct for both peak and
 625 off-peak periods.

626 To enable comparison, we computed the median with error bars using
 627 95% confidence intervals using bootstrapping over all 45 experiments. Fig-
 628 ure 9 shows these results over the supermarket data-set using the k-means
 629 for each one of the representations (whole profile, 8-feat, 4-feat and 2-feat).
 630 The results show only small differences between 2-feat clustering compared
 631 with using the whole profile. Interestingly, for the CDI (Figure 9a) and
 632 SI (Figure 9b) evaluators the clustering 2-feat results outperform those ob-
 633 tained with the whole profile when the number of clusters is greater than

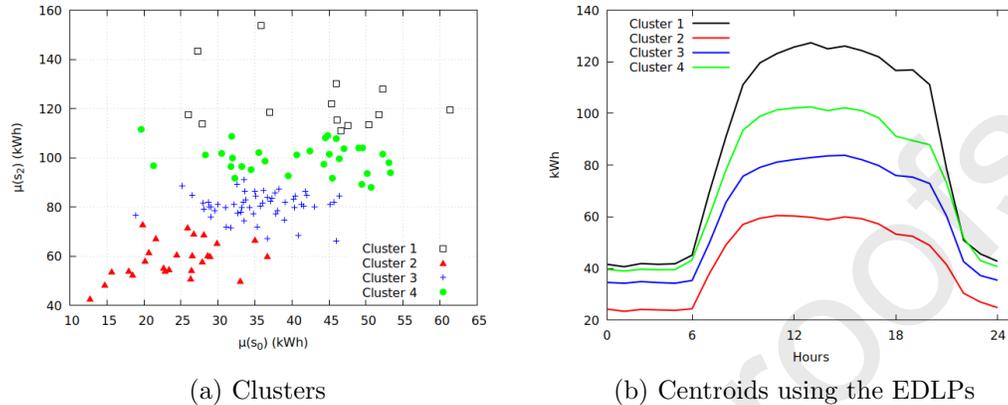


Figure 8: Clustering results for EDLPs represented with $\mu(s_0)$ and $\mu(s_2)$ (only) using data for Winter 2017 of SEG supermarkets with k-means ($k=4$). Clusters 1, 2, 3 and 4 have 15, 26, 57 and 31 points, respectively.

634 three. Generally, 2-feat scores are better than scores obtained with 4-feat
 635 and 8-feat.

636 The fact that the 8-feat results include \vec{t} are worse than the other cluster-
 637 ing results is due to two factors: 1) \vec{t} are numeric variables but they represent
 638 time intervals that are not well modelled by clustering algorithms that use
 639 Euclidean distances, and 2) the time intervals may add noise when creating
 640 the clusters as they are evaluated only using the demand differences of the
 641 whole profile.

642 Clustering results are given in Table 2 for all the evaluators averaged
 643 over all the whole profile (left value) and 2-feat (right value) experiments and
 644 number of cluster separated by algorithm. The differences between the values
 645 are small meaning that the results with both representations are similar. It
 646 might be expected that the whole-profile clustering evaluator would be better
 647 than the 2-feat results, however, for some algorithms and evaluators *e.g.* k-
 648 means and SI, or single link and SI, this is not the case.

649 For the retail store data-set, clustering experiments are performed inde-
 650 pendently for all the EDLPs computed during each season and for the whole
 651 year (Figure 5) with similar characteristics to the supermarkets data-set.
 652 When the number of clusters is small (less than four or five) the differences
 653 between the scores obtained with the whole profile and the reduced feature
 654 representation is greater than when using more clusters. Results obtained

Alg/Eval	CDI	MDI	SI	DBI	MIA	VRC
Kmeans	0.44/0.43	1.26/1.26	24.98/23.97	1.07/1.14	10.59/10.78	134.30/126.02
Single	0.30/0.30	1.14/1.25	6.85/6.64	0.55/0.57	8.88/9.19	10.92/14.47
Complete	0.35/0.36	0.76/0.88	14.89/15.79	0.93/1.00	10.06/10.41	116.63/107.78
UPGMA	0.30/0.31	0.57/0.70	9.29/9.40	0.72/0.83	9.34/9.82	90.40/91.02
WPGMA	0.32/0.33	0.65/0.77	13.57/12.78	0.82/0.90	9.50/9.97	96.47/94.18
UPGMC	0.28/0.29	0.56/0.69	8.99/8.83	0.65/0.78	8.94/9.64	84.64/87.24
WPGMC	0.28/0.30	0.60/0.69	9.82/10.09	0.65/0.80	8.98/9.61	80.44/90.69
WARD	0.47/0.49	1.36/1.51	29.01/29.37	1.14/1.23	10.59/10.98	126.54/118.02

Table 2: Clustering results for the supermarket data-set for all evaluators averaged over all the whole profile (left value), 2-feat (right value), and number of cluster separated by the algorithm.

655 with 8-feat are consistently worse than those obtained with the other rep-
656 resentations. Clustering results are given in Table 3 or all the evaluators
657 averaged over all the whole profile (left value) and 2-feat (right value) exper-
658 iments and number of cluster separated by algorithm. The results obtained
659 with the whole profile marginally outperform those obtained with the 2-feat,
660 with exceptions such a UPGM algorithm and SI evaluator.

Alg/Eval	CDI	MDI	SI	DBI	MIA	VRC
Kmeans	0.21/0.24	2.90/2.56	13.72/17.06	0.87/0.86	1.66/1.77	750.69/734.90
Single	0.09/0.09	0.52/0.66	2.93/2.91	0.21/0.27	0.85/0.91	141.50/145.86
Complete	0.14/0.17	0.72/0.83	3.96/4.71	0.59/0.66	1.35/1.44	497.56/507.80
UPGMA	0.12/0.12	0.44/0.51	3.34/3.38	0.44/0.49	1.19/1.20	278.29/342.59
WPGMA	0.12/0.13	0.54/0.70	3.46/3.49	0.50/0.53	1.19/1.24	370.92/381.44
UPGMC	0.12/0.12	0.45/0.49	3.41/3.13	0.44/0.47	1.16/1.17	308.14/333.41
WPGMC	0.13/0.13	0.51/0.53	3.62/3.64	0.46/0.49	1.24/1.23	270.93/352.74
WARD	0.66/0.79	5.09/7.81	143.51/139.43	1.21/1.31	2.05/2.47	484.12/388.30

Table 3: Clustering results for the retail stores data-set for all evaluators averaged over all the whole profile (left value) and 2-feat (right value) experiments, and number of cluster separated by algorithm.

661 As a final remark about the clustering results, evaluation scores for the 2-
662 feat clustering results are slightly worse than those obtained when using the

663 whole profile when using less than four clusters. However, evaluation scores
664 for these two representations are very close when the number of clusters is
665 greater than four or averaged over the total number of clusters. The 2-feat
666 works well for clustering the profiles because these two features ($\mu(s_0)$ and
667 $\mu(s_2)$) are the main behavioural drivers accounting for most of the EDLP.

668 5. Conclusions and future work

669 Our aim was to investigate whether dimensional reduction could gener-
670 ate a statistically reasonable representation the EDLP of a retail store such
671 that it could be used to predict the electricity demand for a new store in
672 the portfolio of a company. Previously we have shown how this can be done
673 using the whole profile, but a simpler representation of the values of the fea-
674 tures (*e.g.* Figure 4) may offer advantages by reducing the complexity of the
675 problem. In particular, whether it could help detect trends and anomalous
676 behaviours within EDLPs.

677 We have studied the impact to reduced-feature sets to represent EDLPs
678 for prediction and clustering using real data of two distinct data-sets: super-
679 markets with 1-h resolution readings (prediction and clustering) and retail
680 stores with 30-min resolution (clustering only). We have demonstrated that
681 the extracted features give a good description of the original EDLP *i.e.* being
682 able to re-construct the EDLP with only a small error. However, we need
683 to be aware that for a small number of stores (*e.g.* Figure 5c) the proposed
684 representation did not work so well. In general though, we have shown that
685 the evaluation scores are the same or only marginally worse than results ob-
686 tained using the whole profile. The results are robust as the two tasks are
687 different in nature: prediction is supervised learning meanwhile clustering is
688 unsupervised.

689 This proposed simplified representation is a more concise way to represent
690 the EDLP than using the whole EDLP (real resolution values). For some
691 types of analyses, small variances of the demand within the time period can
692 be considered superfluous information that does not add useful information
693 to the overall picture. For example, as the repeated night-time demand values
694 (Figure 1) are repeated over a long period, using an average value is sufficient
695 to summarise and represent the demand during these periods.

696 The main implication for energy managers and researchers is that a re-
697 duced number of features is easier to interpret and visualise instead of a high

698 resolution EDLP. The clustering results suggest its utility as dimensional re-
 699 duction technique to cope with the ‘curse of dimensionality’. More generally,
 700 we have demonstrated that a simpler way to represent data can work as well
 701 for some specific energy problems as complex and high resolution represen-
 702 tation. As modern (networked) sensors increase the volume, availability, and
 703 immediacy, transforming such high-resolution data streams in a ‘smart’ way
 704 based on observed behaviours may be helpful. The proposed features to rep-
 705 resent the EDLP may have limitations for applications such as investigating
 706 and predicting demand shifting and demand variability for energy manage-
 707 ment purposes . This is due to the lack of granularity which will not allow
 708 detection of demand changes at specific times (e.g. hourly).

709 As future work, we suggest that reduced-feature representation can be
 710 applied to any electricity data-set of retail facilities with a diurnal opening
 711 schedule. Moreover, this feature-reduction technique can be applied classi-
 712 fication. Furthermore, combining both clustering and prediction may be an
 713 interesting approach to separately predict the demand by existing buildings
 714 that are in each cluster. This is different to predicting the demand of new
 715 buildings, but large data-sets in both temporal dimensional and number of
 716 stores would be required for such analysis.

717 References

- 718 [1] UNFCCC, COP26 The Glasgow Climate Pact, Glasgow, UK (2021).
 719 URL <https://ukcop26.org/wp-content/uploads/2021/11/COP26-Presidency-Outcomes-The-Climate-Pact.pdf>
 720
- 721 [2] European Commission, A Clean Planet for all: A European strategic long-
 722 term vision for a prosperous, modern, competitive and climate neutral
 723 economy, (No. COM(2018) 773 final). Brussels, Belgium (2018).
 724 URL <https://eur-lex.europa.eu/legal-content/EN/TXT/PDF/?uri=CELEX:52018DC0773&from=EN>
 725
- 726 [3] CCC, Reducing UK emissions: 2019 Progress Report to Parliament.
 727 Committee on Climate Change, London, UK (2019).
 728 URL <https://www.theccc.org.uk/wp-content/uploads/2019/07/CCC-2019-Progress-in-reducing-UK-emissions.pdf>
 729
- 730 [4] J. Rogelj, D. Shindell, K. Jiang, S. Fifita, P. Forster, V. Ginzburg,
 731 C. Handa, H. Kheshgi, S. Kobayashi, E. Kriegler, L. Mundaca, R. Se-

- 732 ferian, M. Vilarino, K. Calvin, O. Edelenbosch, J. Emmerling, S. Fuss,
733 T. Gasser, N. Gillet, C. He, E. Hertwich, L. H. Isaksson, D. Hupp-
734 mann, G. Luderer, A. Markandya, D. McCollum, R. Millar, M. Mein-
735 shausen, A. Popp, J. Pereira, P. Purohit, K. Riahi, A. Ribes, H. Saun-
736 ders, C. Schadel, C. Smith, P. Smith, E. Trutnevyte, Y. Xiu, K. Zickfeld,
737 W. Zhou, Chapter 2: Mitigation pathways compatible with 1.5 °C in the
738 context of sustainable development, in: Global Warming of 1.5 °C an
739 IPCC special report on the impacts of global warming of 1.5 °C above
740 pre-industrial levels and related global greenhouse gas emission path-
741 ways, in the context of strengthening the global response to the threat
742 of climate change, Intergovernmental Panel on Climate Change, 2018,
743 pp. 93–174.
744 URL <http://pure.iiasa.ac.at/id/eprint/15515/>
- 745 [5] E. Bassas, J. Patterson, P. Jones, A review of the evolution of green
746 residential architecture, *Renewable and Sustainable Energy Reviews* 125
747 (2020) 109796.
748 URL <https://doi.org/10.1016/j.rser.2020.109796>
- 749 [6] H. Rau, P. Moran, R. Manton, J. Goggins, Changing energy cultures?
750 Household energy use before and after a building energy efficiency
751 retrofit, *Sustainable Cities and Society* 54 (2020) 101983.
752 URL <https://doi.org/10.1016/j.scs.2019.101983>
- 753 [7] E. Cuce, An overview of domestic energy consumption in the UK: past,
754 present and future, *International Journal of Ambient Energy* (2014) 1–8.
755 URL <https://doi.org/10.1080/01430750.2014.973120>
- 756 [8] R. Granell, C. Axon, D. Wallom, R. Layberry, Power-use profile anal-
757 ysis of non-domestic consumers for electricity tariff switching, *Energy*
758 *Efficiency* 9 (2016) 825–841.
759 URL <https://doi.org/10.1007/s12053-015-9404-9>
- 760 [9] C. Axon, S. Bright, T. Dixon, K. Janda, M. Kolokotroni, Building com-
761 munities: Reducing energy use in tenanted commercial property, *Build-
762 ing Research and Information* 40 (2012) 461–472.
763 URL <https://doi.org/10.1080/09613218.2012.680701>
- 764 [10] F. M. Dahunsi, A. E. Olawumi, D. T. Ale, O. A. Sarumi, A systematic
765 review of data pre-processing methods and unsupervised mining meth-

- 766 ods used in profiling smart meter data, *AIMS Electronics and Electrical*
767 *Engineering* 5 (4) (2021) 284–314.
768 URL <https://doi.org/10.3934/electreng.2021015>
- 769 [11] G. Chicco, Overview and performance assessment of the clustering meth-
770 ods for electrical load pattern grouping, *Energy* Vol. 42 (1) (2012) 68–80.
771 URL <https://doi.org/10.1016/j.energy.2011.12.031>
- 772 [12] S. Yilmaz, J. Chambers, M. Patel, Comparison of clustering approaches
773 for domestic electricity load profile characterisation - implications for
774 demand side management, *Energy* 180 (2019) 665–677.
775 URL <https://doi.org/10.1016/j.energy.2019.05.124>
- 776 [13] R. Granell, C. Axon, M. Kolokotroni, D. Wallom, Predicting electricity
777 demand profiles of new supermarkets using machine learning, *Energy*
778 *and Buildings* 234 (2021) 110635–110635.
779 URL <https://doi.org/10.1016/j.enbuild.2020.110635>
- 780 [14] C. M. Bishop, *Pattern Recognition and Machine Learning*, Springer-
781 Verlag, 2006.
- 782 [15] R. E. Bellman, *Adaptive control processes*, Princeton University Press,
783 2015.
- 784 [16] W. Chung, Review of building energy-use performance benchmarking
785 methodologies, *Applied Energy* 88 (5) (2011) 1470 – 1479.
786 URL <http://doi.org/10.1016/j.apenergy.2010.11.022>
- 787 [17] Z. Li, Y. Han, P. Xu, Methods for benchmarking building energy con-
788 sumption against its past or intended performance: An overview, *Ap-
789 plied Energy* 124 (2014) 325 – 334.
790 URL <http://doi.org/10.1016/j.apenergy.2014.03.020>
- 791 [18] H.-X. Zhao, F. Magouls, A review on the prediction of building energy
792 consumption, *Renewable and Sustainable Energy Reviews* 16 (6) (2012)
793 3586 – 3592.
794 URL <https://doi.org/10.1016/j.rser.2012.02.049>
- 795 [19] C. Deb, F. Zhang, J. Yang, S. E. Lee, K. W. Shah, A review on time se-
796 ries forecasting techniques for building energy consumption, *Renewable*

- 797 and Sustainable Energy Reviews 74 (2017) 902 – 924.
798 URL <https://doi.org/10.1016/j.rser.2017.02.085>
- 799 [20] M. Bourdeau, X.-Q. Zhai, E. Nefzaoui, X. Guo, P. Chatellier, Modeling
800 and forecasting building energy consumption: A review of data-driven
801 techniques, Sustainable Cities and Society 48 (2019) 101533.
802 URL <https://doi.org/10.1016/j.scs.2019.101533>
- 803 [21] K. Li, W. Xue, G. Tan, A. S. Denzer, A state of the art review on the
804 prediction of building energy consumption using data-driven technique
805 and evolutionary algorithms, Building Services Engineering Research
806 and Technology 41 (1) (2020) 108–127.
807 URL <https://doi.org/10.1177/0143624419843647>
- 808 [22] K. Yun, R. Luck, P. J. Mago, H. Cho, Building hourly thermal load
809 prediction using an indexed arx model, Energy and Buildings 54 (2012)
810 225–233.
811 URL <https://doi.org/10.1016/j.enbuild.2012.08.007>
- 812 [23] K. Jeong, C. Koo, T. Hong, An estimation model for determining the an-
813 nual energy cost budget in educational facilities using sarima (seasonal
814 autoregressive integrated moving average) and ann (artificial neural net-
815 work), Energy 71 (2014) 71–79.
816 URL <https://doi.org/10.1016/j.energy.2014.04.027>
- 817 [24] D. W. Schrock, D. E. Clarige, Predicting energy usage in a supermarket,
818 in: Proceedings of the Sixth Symposium on Improving Building Systems
819 in Hot and Humid Climates, 1989, pp. 19–27.
- 820 [25] K. Lindberg, S. Bakker, I. Sartori, Modelling electric and heat load
821 profiles of non-residential buildings for use in long-term aggregate load
822 forecasts, Utilities Policy 58 (2019) 63–88.
823 URL <https://doi.org/10.1016/j.jup.2019.03.004>
- 824 [26] M. S. Spyrou, K. Shanks, M. J. Cook, J. Pitcher, R. Lee, An empirical
825 study of electricity and gas demand drivers in large food retail buildings
826 of a national organisation, Energy and Buildings 68, Part A (2014) 172
827 – 182.
828 URL <http://doi.org/10.1016/j.enbuild.2013.09.015>

- 829 [27] W. Chung, Y. Hui, Y. M. Lam, Benchmarking the energy efficiency of
830 commercial buildings, *Applied Energy* 83 (1) (2006) 1 – 14.
831 URL <http://doi.org/10.1016/j.apenergy.2004.11.003>
- 832 [28] M. Braun, H. Altan, S. Beck, Using regression analysis to predict the
833 future energy consumption of a supermarket in the UK, *Applied Energy*
834 130 (2014) 305 – 313.
835 URL <http://doi.org/10.1016/j.apenergy.2014.05.062>
- 836 [29] O. Valgaev, F. Kupzog, Building power demand forecasting using k-
837 nearest neighbors model - initial approach, in: *2016 IEEE PES Asia-
838 Pacific Power and Energy Engineering Conference (APPEEC)*, 2016,
839 pp. 1055–1060.
- 840 [30] Z. Ma, J. Song, J. Zhang, Energy consumption prediction of air-
841 conditioning systems in buildings by selecting similar days based on
842 combined weights, *Energy and Buildings* 151 (2017) 157 – 166.
843 URL <https://doi.org/10.1016/j.enbuild.2017.06.053>
- 844 [31] M. W. Ahmad, M. Mourshed, Y. Rezgui, Trees vs Neurons: Compar-
845 ison between random forest and ANN for high-resolution prediction of
846 building energy consumption, *Energy and Buildings* 147 (2017) 77 – 89.
847 URL <https://doi.org/10.1016/j.enbuild.2017.04.038>
- 848 [32] G. K. Tso, K. K. Yau, Predicting electricity energy consumption: A
849 comparison of regression analysis, decision tree and neural networks,
850 *Energy* 32 (9) (2007) 1761–1768.
851 URL <https://doi.org/10.1016/j.energy.2006.11.010>
- 852 [33] Y. T. Chae, R. Horesh, Y. Hwang, Y. M. Lee, Artificial neural network
853 model for forecasting sub-hourly electricity usage in commercial build-
854 ings, *Energy and Buildings* 111 (2016) 184–194.
855 URL <https://doi.org/10.1016/j.enbuild.2015.11.045>
- 856 [34] C. Deb, L. S. Eang, J. Yang, M. Santamouris, Forecasting diurnal cool-
857 ing energy load for institutional buildings using artificial neural net-
858 works, *Energy and Buildings* 121 (2016) 284–297.
859 URL <https://doi.org/10.1016/j.enbuild.2015.12.050>
- 860 [35] Q. Li, Q. Meng, J. Cai, H. Yoshino, A. Mochida, Predicting hourly cool-
861 ing load in the building: A comparison of support vector machine and

- 862 different artificial neural networks, *Energy Conversion and Management*
863 50 (1) (2009) 90 – 96.
864 URL <https://doi.org/10.1016/j.enconman.2008.08.033>
- 865 [36] D. Zhao, M. Zhong, X. Zhang, X. Su, Energy consumption predicting
866 model of vrv (variable refrigerant volume) system in office buildings
867 based on data mining, *Energy* 102 (2016) 660 – 668.
868 URL <https://doi.org/10.1016/j.energy.2016.02.134>
- 869 [37] D. Datta, S. Tassou, D. Marriott, Application of neural networks for the
870 prediction of the energy consumption in a supermarket, *Proceedings of*
871 *the Clima 2000 Conference*, Brussels, Belgium.
- 872 [38] G. Revati, J. Hozefa, S. Shadab, A. Sheikh, S. R. Wagh, N. M.
873 Singh, Smart building energy management: Load profile predic-
874 tion using machine learning, in: *2021 29th Mediterranean Con-*
875 *ference on Control and Automation (MED)*, 2021, pp. 380–385.
876 doi:10.1109/MED51440.2021.9480170.
- 877 [39] G. Hafeez, K. S. Alimgeer, I. Khan, Electric load forecasting based on
878 deep learning and optimized by heuristic algorithm in smart grid, *Ap-*
879 *plied Energy* 269 (2020) 114915.
880 URL <https://doi.org/10.1016/j.apenergy.2020.114915>
- 881 [40] X. Luo, L. O. Oyedele, A. O. Ajayi, O. O. Akinade, H. A. Owolabi,
882 A. Ahmed, Feature extraction and genetic algorithm enhanced adaptive
883 deep neural network for energy consumption prediction in buildings,
884 *Renewable and Sustainable Energy Reviews* 131 (2020) 109980.
885 URL <https://doi.org/10.1016/j.rser.2020.109980>
- 886 [41] B. Dong, C. Cao, S. E. Lee, Applying support vector machines to predict
887 building energy consumption in tropical region, *Energy and Buildings*
888 37 (5) (2005) 545 – 553.
889 URL <https://doi.org/10.1016/j.enbuild.2004.09.009>
- 890 [42] S. Paudel, M. Elmitri, S. Couturier, P. H. Nguyen, R. Kamphuis,
891 B. Lacarrere, O. Le Corre, A relevant data selection method for en-
892 ergy consumption prediction of low energy building based on support
893 vector machine, *Energy and Buildings* 138 (2017) 240 – 256.
894 URL <https://doi.org/10.1016/j.enbuild.2016.11.009>

- 895 [43] R. E. Edwards, J. New, L. E. Parker, Predicting future hourly residen-
896 tial electrical consumption: A machine learning case study, *Energy and*
897 *Buildings* 49 (2012) 591–603.
898 URL <https://doi.org/10.1016/j.enbuild.2012.03.010>
- 899 [44] R. K. Jain, K. M. Smith, P. J. Culligan, J. E. Taylor, Forecasting energy
900 consumption of multi-family residential buildings using support vector
901 regression: Investigating the impact of temporal and spatial monitoring
902 granularity on performance accuracy, *Applied Energy* 123 (2014) 168–
903 178.
904 URL <https://doi.org/10.1016/j.apenergy.2014.02.057>
- 905 [45] G. J. Tsekouras, N. D. Hatziargyriou, E. N. Dialynas, Two-stage pattern
906 recognition of load curves for classification of electricity customers, *IEEE*
907 *Transactions on Power Systems* 22 (3) (2007) 1120–1128.
908 URL <https://doi.org/10.1109/TPWRS.2007.901287>
- 909 [46] P. Nystrup, H. Madsen, E. M. Blomgren, G. de Zotti, Clustering com-
910 mercial and industrial load patterns for long-term energy planning,
911 *Smart Energy* 2 (2021) 100010.
912 URL <https://doi.org/10.1016/j.segy.2021.100010>
- 913 [47] G. Chicco, R. Napoli, F. Piglione, Comparisons among clustering tech-
914 niques for electricity customer classification, *IEEE Transactions on*
915 *Power Systems* 21 (2) (2006) 933–940.
916 URL <https://doi.org/10.1109/TPWRS.2006.873122>
- 917 [48] R. Granell, C. J. Axon, D. C. Wallom, Clustering disaggregated load
918 profiles using a Dirichlet process mixture model, *Energy Conversion and*
919 *Management* 92 (2015) 507–516.
920 URL <http://dx.doi.org/10.1016/j.enconman.2014.12.080>
- 921 [49] A. Notaristefano, G. Chicco, F. Piglione, Data size reduction with sym-
922 bolic aggregate approximation for electrical load pattern grouping, *IET*
923 *Generation, Transmission Distribution* 7 (2) (2013) 108–117.
- 924 [50] R. Granell, C. J. Axon, D. C. H. Wallom, Impacts of raw data temporal
925 resolution using selected clustering methods on residential electricity
926 load profiles, *IEEE Transactions on Power Systems* 30 (6) (2015) 3217–
927 3224.
928 URL <https://doi.org/10.1109/TPWRS.2014.2377213>

- 929 [51] S. Roberts, J. Thunim, Managing and mining smart meter data - at
930 scale, cSE Project Showcase Presentation. Centre for Sustainable Energy
931 (2013).
- 932 [52] Z. Mylona, M. Kolokotroni, S. A. Tassou, Frozen food retail: Measuring
933 and modelling energy use and space environmental systems in an oper-
934 ational supermarket, *Energy and Buildings* 144 (2017) 129 – 143.
935 URL <http://doi.org/10.1016/j.enbuild.2017.03.049>
- 936 [53] C. Cortes, V. Vapnik, Support-vector networks, in: *Machine Learning*,
937 1995, pp. 273–297.
- 938 [54] F. Hayashi, *Econometrics*, Princeton University Press, 2011.
- 939 [55] S. Fritsch, F. Guenther, *Neuralnet: Training of Neural Networks*, R
940 package version 1.33 (2016).
941 URL <https://CRAN.R-project.org/package=neuralnet>
- 942 [56] D. Meyer, E. Dimitriadou, K. Hornik, A. Weingessel, F. Leisch, e1071:
943 Misc Functions of the Department of Statistics, Probability Theory
944 Group (Formerly: E1071), TU Wien, r package version 1.6-8 (2017).
945 URL <https://CRAN.R-project.org/package=e1071>
- 946 [57] J. Williams, Clustering household electricity use profiles, in: *Proceed-*
947 *ings of MLSDA*, ACM, 2013, pp. 19–26.
- 948 [58] T. Räsänen, D. Voukantsis, H. Niska, K. Karatzas, M. Kolehmainen,
949 Data-based method for creating electricity use load profiles using large
950 amount of customer-specific hourly measured electricity use data, *Ap-*
951 *plied Energy* 87 (11) (2010) 3538–3545.
952 URL <https://doi.org/10.1016/j.apenergy.2010.05.015>
- 953 [59] M. Halkidi, Y. Batistakis, M. Vazirgiannis, On clustering validation
954 techniques, *J. Intell. Inf. Syst.* 17 (2-3) (2001) 107–145.

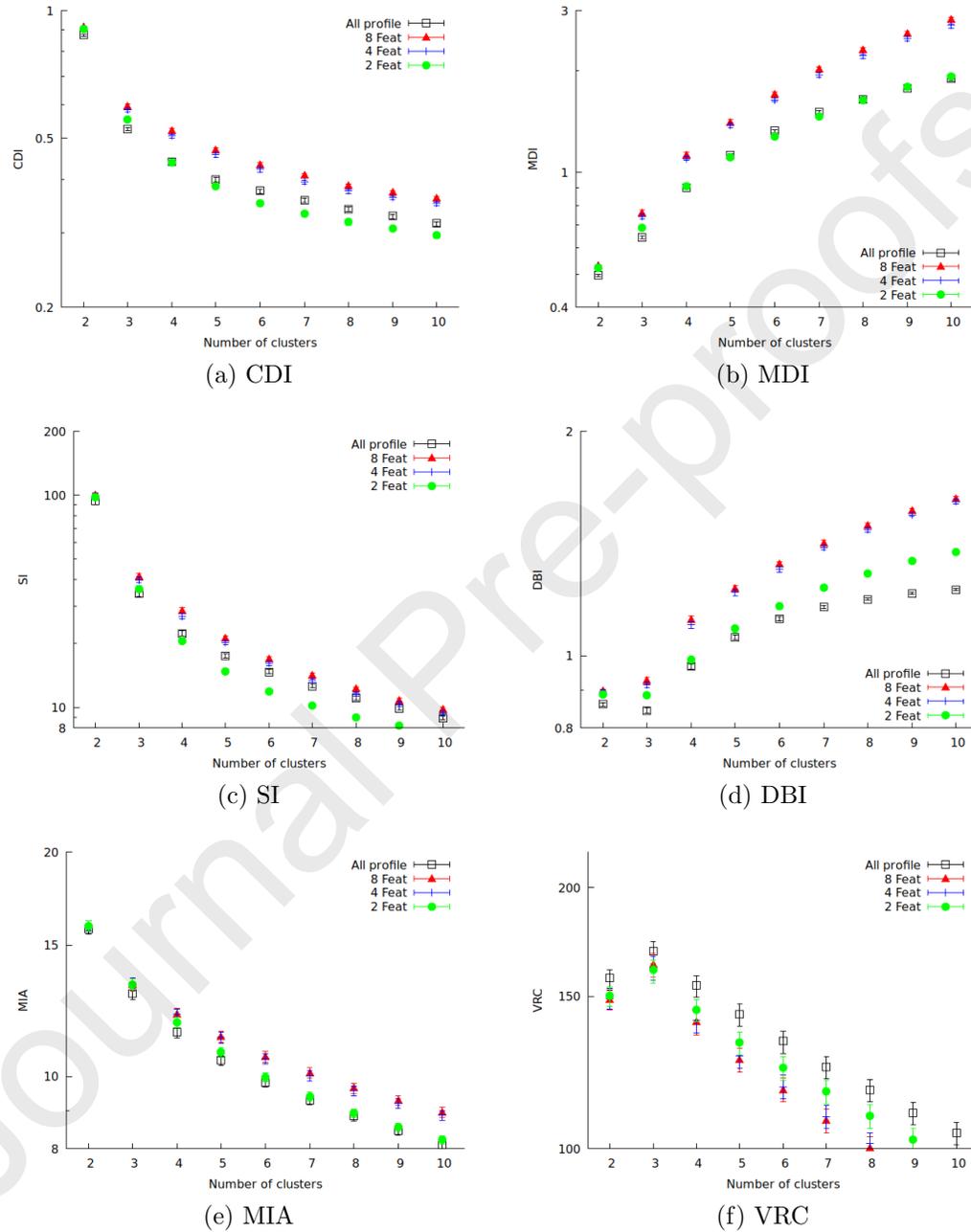


Figure 9: Clustering results for the supermarket data-set using the k-means. N.B. the Y-axis is log scale.

Declaration of interests

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

The authors declare the following financial interests/personal relationships which may be considered as potential competing interests:

Journal Pre-proofs