

Classification of new electricity customers based on surveys and smart metering data

Joaquim L. Viegas^{a,*}, Susana M. Vieira^a, R. Melício^{a,b}, V. M. F. Mendes^{b,c},
João M. C. Sousa^a

^a*IDMEC, LAETA, Instituto Superior Técnico, Universidade de Lisboa,
Av. Rovisco Pais, 1, 1049-001 Lisbon, Portugal*

^b*Dep.de Física, Escola de Ciências e Tecnologia, Universidade de Évora, Portugal*

^c*Instituto Superior de Engenharia de Lisboa*

Abstract

This paper proposes a process for the classification of new residential electricity customers. The current state of the art is extended by using a combination of smart metering and survey data and by using model-based feature selection for the classification task. Firstly, the normalized representative consumption profiles of the population are derived through the clustering of data from households. Secondly, new customers are classified using survey data and a limited amount of smart metering data. Thirdly, regression analysis and model-based feature selection results explain the importance of the variables and which are the drivers of different consumption profiles, enabling the extraction of appropriate models. The results of a case study show that the use of survey data significantly increases accuracy of the classification task (up to 20%). Considering four consumption groups, more than half of the customers are correctly classified with only one week of metering data, with more weeks the accuracy is significantly improved. The use of model-based feature selection resulted in the use of a significantly lower number of features allowing an easy interpretation of the derived models.

Keywords:

Data-driven energy efficiency, Electricity customer clustering, Classification of new residential customers, Customer feature selection, Smart metering data, Customer surveys data

*Corresponding author.

Email address: joaquim.viegas@tecnico.ulisboa.pt (Joaquim L. Viegas)

1. Introduction

A game-changing shift has been happening in the utility industry and energy markets. Policy focused on energy efficiency and sustainability is growing fruit of the awareness of current environmental challenges. Liberalization, growing competition between utilities, technological advancements and policy towards a sustainable use of energy sources are pushing utilities to seek innovation and new market related insights.

Electricity is a main energy carrier used around the world for supporting the primary, secondary and tertiary sectors. The commercial and residential energy demand is expected to continue to shift towards electricity and away from primary fuels. By 2040, forecasts indicate that electricity generation will account for more than 40% of global energy consumption and, from 2010 to 2040, global electricity demand is projected to increase by about 85% [1–3].

Technological advancement in the fields of metering, communications and computation are enabling utilities to monitor and save huge amounts of data related to their operation. The deployment of electricity meters with two-way communication capabilities is enabling the logging of the consumption of users with high resolution. The number of advanced metering infrastructure (AMI) installations, also known as smart meters, has surpassed the number of traditional one-way communication meters in the United States [4]. Close to 45 million smart meters are already installed in three Member States (Finland, Italy and Sweden) of the European Union (EU), representing 23 percent of the envisaged installation in the EU by 2020 [5].

The consumption data of customers has the potential to give insights of great importance for utilities and policy makers. Valuable insights can be derived by the knowledge of typical consumption curves of different consumer groups and understanding what are the main drivers of consumption. This knowledge can assist decision makers in the electricity utility industry in developing demand side management (DSM) programs, consumer engagement strategy, marketing, alternative tariff setting methods and demand forecasting tools [6]. Knowledge on the way different demographic groups consume electricity is valuable to study the effect of energy policy on different population groups.

The high number of consumers and desired high sampling frequencies in smart metering implies that huge amounts of data have to be stored and processing grows in complexity. [Computational intelligence techniques in the fields of machine learning are starting to be extensively used in order](#)

to extract knowledge from the data coming from the grid. These techniques can provide decision makers with predictive models and the ability to extract valuable knowledge.

In order to characterize the behaviour of electricity customers, the clustering of electricity consumption data has been the focus of a considerable amount of research in the past years. The usual stated applications range from the design and simulation of DSM [7, 8], load forecasting [9–11], tariff setting [12–14], marketing and bad data detection. The clustering methods found to be used are mostly the K-means algorithm [8, 15–18]. Fuzzy clustering [19] has shown promise in the field. Data preparation is of high importance in these applications, dictating what information is desired to be extracted from the clustering and the ability of the used methods to achieve good results. Normalization, parametric modelling [10], temperature based normalization [16, 20] and wavelet transformation [9] have been found to be used in the literature.

The use of static data related to household characteristics, e.g., income, number of inhabitants, education, construction year and appliances in relation to static or dynamic energy consumption data is being studied in order to find the main drivers of residential energy consumption. In [21–23] factor analysis and linear regression are used to find the main determinants of energy consumption in residential settings, such as weather data, household characteristics and demographics. In [24] demographic data and psychological and belief related data is studied in comparison to energy consumption. [25, 26] presents studies on the prediction of household information based on smart meter data. In [27, 28] consumptions profiles obtained via clustering are correlated to household characteristics. In [29] a methodology is presented for the characterization of medium voltage electricity customers through clustering and posterior modelling for which the classification of new customers is stated as a possible application.

Classifying new customers is crucial for marketing purposes, as customers with lengthy relationships are less likely to defect and are less affected by new information and offers. Thus, a greater impact of marketing strategies and engagement is expected with new customers [30, 31].

This paper extends the current state of the art by developing a process for the classification of new electricity customers using not only metering data but also using static data on household characteristics. The use of a limited amount of metering data is done in order to emulate the analysis of new electricity customers for which only a small amount of data is available.

The use of model-based feature selection for the discovery of the consumption drivers shows promise in the field.

Based on the clustering of customers' electricity consumption data, the consumption profile of new customers is predicted using survey data and a limited amount of smart metering data. Classification models in combination with model-based and filter feature selection are compared for the classification task, selection and analysis of variables.

The developed process aims to provide an interpretable classification modelling method for the classification of electricity customers and discovery of the drivers of different electricity consumption profiles. The presented results aim to illustrate the application of the proposed process, using data that resulted from smart metering trials encompassing more than three thousand households in Ireland [32]. Requirements for the classification of customers and insights on the drivers of residential electricity consumption are presented.

This paper is organized as follows: Section 2 discusses the uses of the proposed process in the context of the smart grid. Section 3 presents the method for the generation of the populations representative consumption profiles. Section 4 presents the techniques used for modelling, feature selection and model evaluation. Section 5 presents the experimental results and presents the discussion and Section 6 presents the conclusions.

2. Classification of customers in the smart grid

The smart grid is a concept with the purpose of intelligently integrating the generation, transmission and consumption of electricity through technological means [33–37]. A smart electricity grid enables an efficient management of the whole electricity supply chain through innovative applications. The applications can provide the capacity to: securely integrate more renewable energy sources and distributed generation; deliver power in a more efficient and secure manner through advanced control and monitoring; automatically reconfigure the grid to prevent and restore outages; better integrate consumption through DSM; enable consumer engagement in the market [38–41].

Smart metering roll-outs and pilots are paving the way for the development of the smart grid. Meters with two-way communication capabilities are expected to empower consumers by enabling the creation of consumer services and engaging them to actively participate in the electricity market.

112 In Europe the total investment of smart grids amounted to €3.15 billion in
113 2014 and smart metering projects account for most of the total investment
114 [38].

115 The imperative for consumers to be on board is defended in order not
116 only to reap the benefits of a smart grid, but also to make smart metering
117 projects profitable. The extent of the transformation of the grid rests on
118 the needs and the willingness of consumers to pay for the implementation
119 [38, 41]. The right consumers need to be identified, engaged and motivated
120 in order to reap the benefits of smart metering in terms of electricity cost
121 savings, through, e.g., load shifting [42].

122 Knowledge on the ways electricity is consumed in a population and what
123 are the drivers of consumption dynamics, e.g., demographics, household char-
124 acteristics and the use of appliances is essential in order to personalize ap-
125 plications, energy services and policy towards a smarter grid.

126 In the context of the smart grid, the ability to effectively group customers
127 into similar behaviour market segments and to find the segment of new cus-
128 tomers is very valuable, e.g., in the following applications:

- 129 • Proposing tariff offers or DSM schemes taking into account the expected
130 consumption behaviour of the customers;
- 131 • Planning and studying the potential impact of personalized services
132 and offers;
- 133 • Offering the energy saving and sustainability services the customers are
134 most likely to be interested in.

135 The proposed process for clustering and classification of electricity cus-
136 tomers enables more effective customer engagement on the part of utilities
137 and smart grid operators. Customer engagement is essential to maximize the
138 willingness of customers to pay for the implementation of this type of grid,
139 either directly or indirectly by increasing the grids efficiency through DSM
140 programs and energy efficiency solutions.

141 3. Clustering

142 Clustering methods attempt to group objects based on a definition of
143 similarity. The objective is to find groups of objects with greater similarity
144 between them than to the objects of other groups.

145 In the scope of this paper and the analysis of customers' representa-
 146 tive consumption profiles, clustering methods are used to find which are the
 147 groups of customers which have similar consumption curves in some context,
 148 e.g., season, type of day. These groups are represented by the populations
 149 representative consumption profiles, resulting from aggregating the profile of
 150 all the customers of a group, equivalent to the cluster centroid.

151 The methodology followed to find the customer groups and respective
 152 representative consumption profiles is in Figure 1. The clustering process
 153 is similar to the one proposed in [29]. Firstly, smart metering data is pro-
 154 cessed in order to obtain the customers' representative consumption profiles,
 155 secondly, various clustering configurations are tested. Configurations are
 156 evaluated using multiple clustering validity indexes (CVI) which are used,
 157 together with careful visual evaluation, to chose the final configuration and
 158 obtain the customer groups and profiles.

159 3.1. Customers' normalized representative consumption profiles

160 Smart metering consumption data is composed of a large set of times-
 161 tamped intervals with consumption values. In order to obtain consumption
 162 profiles which can be easily interpreted, visualized and manipulated, the data
 163 goes through a process of context filtering, aggregation and pre-processing.

164 The process of context filtering consists on selecting data which represents
 165 a specific context, defined, for example, by a temporal window (e.g. Winter,
 166 Summer), type of day (e.g. working day) and location.

167 Let \mathbf{x}_i be the feature vector (list of variables) associated to customer i .
 168 $\mathbf{x}_i = (\mathbf{x}_i^m, \mathbf{x}_i^s)$ where \mathbf{x}_i^m has dimension r equal to the number of variables
 169 which characterize a customers representative load profile (LP) or derived
 170 load indices (LI) and \mathbf{x}_i^s has dimension t equal to the number of survey vari-
 171 ables used. The dimension of a customers feature vector \mathbf{x}_i is $p = r + t$. The LI
 172 and survey variables are presented in 5.1 and 5.3. $X = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\} \subseteq \mathbb{R}^p$
 173 is the feature dataset of N customers.

174 After filtering, the consumption data is aggregated in order to reduce the
 175 dimension and obtain a curve representative of the whole temporal window.
 176 The aggregation is characterized by the period used, e.g., hourly, daily and
 177 operator, e.g., mean, median. For example, doing an hourly mean aggrega-
 178 tion of the consumption data of customer i will generate a vector $\mathbf{x}_i^m \in \mathbb{R}^{24}$
 179 in which each element represents the mean consumption in a certain hour.

180 The final pre-processing consists on the normalization of the data for eas-
 181 ier clustering, modelling and representation of different information. This

182 paper focuses on the case of normalization for each customer in which each
 183 representative profile is normalized with the maximum value of the profile
 184 as normalization factor. The normalization is done with the intent of trans-
 185 lating the consumption dynamic in relation to the maximum. This is done
 186 in [27–29]. The clustering of absolute representative consumption profiles
 187 results, using the same kind of data, on a separation of groups by amount of
 188 consumption. Without normalization the different shapes of curves are seem-
 189 ingly overshadowed by the mean absolute consumption while clustering [43].

190 Figure 2 pictures an example of the clustering results, showing clusters
 191 centroids for hourly aggregated absolute and normalized representative pro-
 192 files. The curves behave in a similar way for different scales in absolute pro-
 193 files. For normalized consumption profiles the curves are distinct in terms of
 194 linearity and consumption between different times of the day.

195 3.2. *K-means clustering*

196 The K-means algorithm [44] is used due to its simplicity, efficiency and
 197 scalability. The algorithm has been proven to be adequate for this type of
 198 application in the literature [8, 15–18, 45, 46]. Let $\mathbf{S} = \{S_1, \dots, S_J\}$ be the
 199 groups (sets) of customers clustered together, J the number of clusters and
 200 d_e a chosen distance measure. The centroid of a cluster S_k is its mean vector,
 201 $\mu_k = \frac{1}{|S_k|} \sum_{\mathbf{x} \in S_k} \mathbf{x}$. The algorithm is an iterative refinement method which, in
 202 this application, minimizes the distance between the customers’ consumption
 203 profiles \mathbf{x} and the populations μ_k , as given by (1).

$$\arg \min_{\mathbf{S}} \sum_{k=1}^J \sum_{\mathbf{x} \in S_k} d_e(\mathbf{x}, \mu_k)^2 \quad (1)$$

204 The difficulty associated with this algorithm is the need to determine the
 205 number of clusters and their initial centres. The choice of the number of
 206 cluster centres is detailed in the following Section 3.3. The initial cluster
 207 centres are generated randomly and the best clustering result of an high
 208 number of runs is used.

209 3.3. *Clustering evaluation*

210 A clustering in X is a set of disjoint clusters that partition X into k
 211 groups: \mathbf{S} where $\cup_{S_k \in \mathbf{S}} S_k = X, S_k \cap S_l = \emptyset \forall k \neq l$. The euclidean distance
 212 is used and $d_e(\mathbf{x}_i, \mathbf{x}_k) = \sqrt{\sum_{j=1}^p (x_{ij} - x_{kj})^2}$.

213 As pictured in Figure 1, multiple CVI are used to evaluate a number
 214 of different clustering configurations. If there is no consensus between the
 215 different CVI the expert chooses the best configuration based on the analysis
 216 of the CVI and visualization of the clustering results.

217 Three different CVI are used in this work, they evaluate the goodness
 218 of the clustering in terms of maximization of inter cluster distances and
 219 minimization of intra cluster distances [47].

220 The Dunn index (D) [48] is a ratio-type index where the cohesion is esti-
 221 mated by the nearest neighbour distance and the separation by the maximum
 222 cluster diameter. The original index is defined as,

$$D(\mathbf{S}) = \frac{\min_{S_k \in \mathbf{S}} \{\min_{S_l \in \mathbf{S} \setminus S_k} \{\delta(S_k, S_l)\}\}}{\max_{S_k \in \mathbf{S}} \{\Delta(S_k)\}} \quad (2)$$

223 where,

$$\delta(S_k, S_l) = \min_{\mathbf{x}_i \in S_k} \min_{\mathbf{x}_j \in S_l} \{d_e(\mathbf{x}_i, \mathbf{x}_j)\} \quad (3)$$

$$\Delta(S_k) = \max_{\mathbf{x}_i, \mathbf{x}_j \in S_k} \{d_e(\mathbf{x}_i, \mathbf{x}_j)\}. \quad (4)$$

224 The Davis-Bouldin index (DB) [49] estimates the cohesion based on the
 225 distance from the points in a cluster to the centroid and the separation based
 226 on the distance between centroids. The DB index is defined as:

$$DB(\mathbf{S}) = \frac{1}{J} \sum_{S_k \in \mathbf{S}} \max_{S_l \in \mathbf{S} \setminus S_k} \left\{ \frac{F(S_k) + F(S_l)}{d_e(\mu_k, \mu_l)} \right\} \quad (5)$$

227 where,

$$F(S_k) = \frac{1}{|S_k|} \sum_{\mathbf{x}_i \in S_k} d_e(\mathbf{x}_i, \mu_k). \quad (6)$$

228 The silhouette index (Sil) [50] is a normalized summation-type index.
 229 The cohesion is measured based on the distance between all the points in the
 230 same cluster and the separation is based on the nearest neighbor distance.
 231 The silhouette index is defined as:

$$Sil(\mathbf{S}) = \frac{1}{N} \sum_{S_k \in \mathbf{S}} \sum_{\mathbf{x}_i \in S_k} \frac{b(\mathbf{x}_i, S_k) - a(\mathbf{x}_i, S_k)}{\max\{a(\mathbf{x}_i, S_k), b(\mathbf{x}_i, S_k)\}} \quad (7)$$

232 where,

$$a(\mathbf{x}_i, S_k) = \frac{1}{|S_k|} \sum_{\mathbf{x}_j \in S_k} d_e(\mathbf{x}_i, \mathbf{x}_j) \quad (8)$$

$$b(\mathbf{x}_i, S_k) = \min_{S_l \in \mathbf{S} \setminus S_k} \left\{ \frac{1}{|S_l|} \sum_{\mathbf{x}_j \in S_l} d_e(\mathbf{x}_i, \mathbf{x}_j) \right\}. \quad (9)$$

233 4. Modelling

234 4.1. Classification

235 This work intends to train models to predict the group of a new customer,
236 characterized by a representative consumption profile. Figure 3 pictures the
237 electricity customer classifier.

238 Features are extracted from the survey responses and smart metering data
239 of the customer. Based on the features the classifier returns a categorical
240 variable y indicative of the customer group in which the customer best fits.

241 The classifier is a function φ which maps the features of a customer to
242 a categorical variable y , representing one of the J customer groups. It is
243 defined as:

$$\varphi : \mathbb{R}^p \mapsto y \quad (10)$$

$$y \in \{c_1, c_2, \dots, c_J\} \quad (11)$$

244 Classifiers are trained using the group labels extracted through the clus-
245 tering of a full year of smart metering data, considered as the ground truth to
246 be inferred from features extracted from a limited amount of smart metering
247 data and survey data.

248 The two following sections present the modelling approaches used in this
249 methodology.

250 4.1.1. Logistic regression

251 The logistic regression (LR) models the posterior probabilities of the
252 J classes via linear function in x while ensuring the sum to one and re-
253 maining in $[0, 1]$. The LR model has the form presented in (12), where
254 D represents the input vector [51, 52]. The parameter set of the model is
255 $\theta = \{\beta_{10}, \beta_1^T, \dots, \beta_{(J-1)0}, \beta_{(J-1)}^T\}$.

$$\begin{aligned}
\log \frac{\Pr(y = 1|D = \mathbf{x})}{\Pr(y = J|D = \mathbf{x})} &= \beta_{10} + \beta_1^T \mathbf{x} \\
\log \frac{\Pr(y = 2|D = \mathbf{x})}{\Pr(y = J|D = \mathbf{x})} &= \beta_{20} + \beta_2^T \mathbf{x} \\
&\vdots \\
\log \frac{\Pr(y = J - 1|D = \mathbf{x})}{\Pr(y = J|D = \mathbf{x})} &= \beta_{(J-1)0} + \beta_{(J-1)}^T \mathbf{x}
\end{aligned} \tag{12}$$

256 Using the LR model, if the clustering analysis results in J customer
257 groups, the classifier linearly separates each one of $J - 1$ customer groups to
258 the J customer group.

259 LR is usually fit by maximum likelihood, in the case of the results pre-
260 sented in this paper the Newton-Raphson optimization method is used. For
261 the case of two classes the parameters of the model can be easily interpreted
262 through the significance and sign. In the case of multiple classes the inter-
263 pretation of the model parameters is more complex due to a total set of $J - 1$
264 parameters for each variable.

265 The LR model is chosen due to the simplicity (explained by linear func-
266 tions) and interpretability, enabling the understanding of the role of the dif-
267 ferent input variables in explaining the outcome [51]. Models with increased
268 complexity, such as artificial neural networks or support vector machines,
269 may provide higher accuracy but lack the transparency of the LR model
270 [53].

271 4.1.2. *Decision trees*

272 Binary decision tree (DT) learning consists on fitting data to a tree-like
273 structure. This type of method partitions the feature space into a set of
274 rectangles and usually fits a constant in each one. This paper makes use
275 of the popular tree-based regression and classification method called CART
276 (Classification And Regression Tree) [51]. Tree-based methods have the ad-
277 vantage of an easy interpretation and can be transformed into a simple set
278 of rules if the number of branches is low.

279 In order to grow a classification DT the learning algorithm automatically
280 splits the data into two sets at each level, optimizing some criterion which
281 translates the model accuracy. In this paper the Gini index is used, which is
282 a measure of how often a randomly chosen element from the set is incorrectly

labelled if it is randomly labelled according to the distribution of labels in the subset. The learning algorithm minimizes the difference of this measure between tree levels through the growth of the DT. Using DT in the multiple class case is straightforward and each end node of the tree will give a probability for the J labels. Figure 4 pictures an example of a partition obtained by binary splitting and corresponding DT.

A classification DT model is chosen, similarly to the LR model, due to its interpretability, providing a popular binary tree representation [51].

4.2. Feature selection

The objective of feature selection (FS) is to choose a subset of the available features by eliminating features with little or no predictive information and also redundant features that are strongly correlated [54]. FS techniques are usually divided into filter, wrapper and embedded methods. Wrapper and embedded are usually referred to as model-based methods and filter techniques as model-free methods.

Filter techniques assess the relevance of features by looking only at the intrinsic properties of the data. [Filter techniques are normally easily scalable to very high-dimension datasets and computationally simple, having the disadvantage of not taking into account the interaction with the classifier \[55\].](#)

Wrapper methods embed the classification model within the feature subset search. The selected set of features is obtained by training and testing a specific classification model, rendering this approach tailored to a specific classification algorithm [55].

4.2.1. Regression based filter feature selection

In regression analysis parameters are determined indicating the relationship between the features and the model output. The p -values of the hypothesis tests based on the parameters' standard errors indicate if the corresponding variables are believed to be significantly different from 0 (rejected null hypothesis), thus indicators of the output variable. The regression feature selection method used removes the variables for which the corresponding parameters result in a p -value higher than a certain significance level (5%).

This parametric filter FS technique has been used in multiple studies, together with LR or probit regression, in order to find which are the features which are indicative of a specific electricity consumption profile and are determinants of electricity consumption [22, 23, 28].

318 4.2.2. *Wrapper feature selection*

319 This paper proposes the use of greedy wrapper FS methods to find rela-
320 tions between the characteristics of customers and the typical consumption
321 profile. FS is also done in order to generate interpretable models by signifi-
322 cantly reducing the number of features used to classify new customers.

323 Sequential forward selection and sequential backward elimination [56] are
324 the FS methods used. The forward FS algorithm sequentially selects features,
325 starting with a empty set, choosing the features that improve the most the
326 prediction accuracy. This is done until there is no more improvement in
327 prediction. The backward FS algorithm starts with the full set of features and
328 sequentially removes the ones which result in an improvement in prediction
329 accuracy.

330 4.3. *Model evaluation*

331 In order to maximize the significance of the performance results of the
332 trained classifiers k -fold cross-validation is used [51, 53]. This model vali-
333 dation technique randomly divides the dataset into k folds. The classifier
334 is then trained (using $k - 1$ folds) and evaluated (using 1 fold) k times, as
335 pictured in Figure 5. The modelling approach is then evaluated through the
336 mean and standard deviation of the accuracy.

337 In order to do an unbiased FS the methods presented in Section 4.2
338 are used only based on the training sets so that the process is totally in-
339 dependent from the test data. The wrapper FS methods also make use of
340 cross-validation to evaluate the feature subsets.

341 5. Results and discussion

342 5.1. *Dataset*

343 The proposed methodology is applied to data from 4232 Irish households
344 monitored for one and a half year. The dataset consists of electricity con-
345 sumption data logged at 30 minute intervals and surveys responded before
346 the start of the trial. This dataset resulted from an electricity customer be-
347 haviour trial by the Irish Commission for Energy Regulation (CER). The data
348 is stored and maintained by the Irish Social Science Data Archive (ISSDA)
349 [32].

350 The mean hourly consumption for the four seasons is pictured in Figure
351 6. Consumption follows the typical residential dynamic with a small peak in
352 the morning and lunch time, a larger one at the end of the afternoon and

low consumption during the night. As expected, the mean consumption in winter presents the highest values due to the heating needs.

The distribution of the survey responses on social class and number of children per household is pictured in Figure 7. AB is upper middle class and middle class, C1 is lower middle class, C2 is skilled working class, DE is working and non-working classes and F represents farmers. The distributions show that the used data encompasses different demographic groups and household types.

The survey questions used as features are presented in Table 1 to Table 4, along with a description and possible responses. Table 1 presents the features with information on the respondent, Table 2 is related to the habitation characteristics, Table 3 to the heating systems and Table 4 to the appliances.

Survey variables with no response are considered as 'refused'. The customers not considered in the study are the ones who did not respond to the question indicating the number of adults in the household. The final dataset used contains 3440 electricity customers.

5.2. Clustering

This section presents the results from the extraction of features from the customers smart metering data, transformation in representative profiles and clustering in order to obtain the final populations representative consumption profiles.

5.2.1. Customers' representative consumption profiles

In order to obtain the customers' consumption profiles the parameters used to extract the representative features are:

- Context: Only the smart metering data from working days is used and profiles are extracted seasonally;
- Aggregation: The data is aggregated hourly resulting in twenty-four features ($r = 24$);
- Operator: The operator used is the mean.
- Normalization: The profiles are normalized with regards to each customers maximum hourly consumption.

The final customers' representative consumption profiles are equal to the customer normalized mean hourly consumption in working days. The profiles are obtained for each one of the four seasons.

387 5.2.2. Populations representative consumption profiles

388 Following the proposed methodology, the best number of clusters is found
389 to be equal to four for the four seasons. Figure 8 pictures the evolution of
390 the three CVI used when generating between two and six clusters for the
391 Winter season. The Silhouette, Dunn and Davis-Bouldin indexes indicate,
392 respectively, that the best number of cluster is two, four and five. In order to
393 choose a number of clusters the partitions are visually analysed as pictured
394 in Figure 9, Figure 10 and Figure 11. The figures present the populations
395 representative consumption profiles (cluster centres) and the customers' rep-
396 resentative consumption profiles pertaining to the cluster.

397 With two clusters, as pictured in Figure 9, many customers have a con-
398 sumption profile different from the centre, indicating the need for an higher
399 number of clusters. With four clusters, as pictured in Figure 10, the clusters
400 are sufficiently compact having a significant number of customers in each
401 group. With five clusters, as pictured in Figure 11, *Cluster 2* has a low num-
402 ber of customers with profiles showing a low similarity. Based on the visual
403 analysis the number of chosen clusters is equal to four. The same process is
404 used for the other seasons.

405 The final populations representative consumption profiles are pictured in
406 Figure 12. The population is divided mainly due to the following consump-
407 tion profile characteristics:

- 408 • **Peakiness:** Relation between peak evening consumption and the con-
409 sumption throughout the rest of the day. For example: in Winter,
410 clusters 1 and 2 have a much higher difference between peak evening
411 and the rest of the days consumption (high *peakiness*), in comparison
412 to clusters 3 and 4 (low *peakiness*).
- 413 • **Decline time:** Time at which the consumption starts to rapidly de-
414 cline after peak evening consumption. For example: in Spring, clusters
415 2 and 4 have a late declining consumption (late decline) in comparison
416 to clusters 1 and 3 (early decline), specially cluster 3 that has a very
417 early decline in consumption.
- 418 • **Off-peak consumption:** Presence of significant consumption during
419 the off-peak hours (night and early morning) in comparison to the rest
420 of the day. For example: in Autumn, cluster 4 presents a significant
421 consumption during the night hours (high off-peak consumption) in
422 comparison to the clusters 1, 2 and 3 (low off-peak consumption).

Summer presents the most different populations consumption profiles in comparison to the other seasons, as pictured by the the consumption profile of *Cluster 2*. This cluster presents a high amount of variability between customers results in a low mean normalized consumption throughout the day.

Table 5 presents the distribution of customers between the different clusters for each one of the seasons. Asides from the Winter clustering, the customers are approximately uniformly distributed between the four groups.

5.3. Classification of new customers and feature selection

Features extracted from metering data and from conducted surveys are used for the classification of new customers. In order to evaluate the process for the classification of new customers, the metering data is limited to an amount starting from no data to ten weeks of data. Due to the high amount of metering data and desire for interpretable models two types of features extracted from the smart metering data are tested: load profile (LP) and load indices (LI).

The LP features are the ones used in the clustering: in this paper they are the hourly aggregated mean consumption normalized on an individual basis. The features differ from the ones used for clustering due to being derived from a limited amount of smart metering data.

The LI are shape indices derived from the LP, these are proposed in [57] and used for the characterization of medium-voltage customers in [29]. LI are used in this paper with the intention of obtaining models of easier interpretation, explaining what consumption characteristics are the most relevant when comparing customers. The indices are presented in Table 6. i_1 is the load factor, i_2 is the off-peak factor, i_3 is the night impact coefficient, i_4 is the lunch impact coefficient and i_5 is the modulation coefficient at off-peak hours. P_{max} , P_{min} , P_{av} are, respectively, the maximum, minimum and average consumption of the corresponding periods.

Table 7 summarizes the smart metering features used in classification. In the case at least one day of metering data is available, a total of $p = r + t = 24 + 47 = 71$ features are available using the LP as the smart metering features and $p = 5 + 47 = 52$ features are available using the LI.

Table 8 and Table 9 present the mean and standard deviation of the accuracy of the trained classifiers, through 5-fold cross-validation, in the cases of no smart metering data, 1, 4, 8 and 10 weeks of available smart metering data (W). In parentheses the mean number of features selected is

460 presented. The results are presented for the LR and DT models, for each
461 season, and further divided by the use of no FS, the filter FS algorithm and
462 forward FS. Backward FS results in a performance closely similar to the use
463 of no FS. Accuracy was used, instead of measures that can correctly deal
464 with class imbalances, such as the Area Under the ROC Curve (AUC) [58],
465 precision/recall and MCC, due to the multiclass nature of the classification
466 problem and the approximately balanced nature of the classes, inferred from
467 Table 5.

468 The evolution of the LR classifier performance with a growing number of
469 weeks of metering data for the Winter season is pictured in Figure 13. The
470 figure shows that, when using LP, the classification accuracy always benefits
471 from the use of survey features. The difference between the performance
472 of the classifier with and without survey features grows with the number
473 of available weeks of smart metering data. When using LI the difference is
474 only significant for the case when there is not metering data for which the
475 classification is random because no features are available.

476 Based on the analysis of the results of Table 8 and Table 9, the use
477 of LP results in an better classification performance, proving that the LI
478 are not able to correctly translate all the information needed to classify the
479 customers.

480 In general, filter FS results in the best accuracy, reducing significantly
481 the number of features in comparison with not using any FS. Using forward
482 FS resulted in an even greater reduction of the number of features at the cost
483 of a reduction of accuracy.

484 The following paragraphs present a detailed analysis of the classification
485 and feature selection results for:

- 486 1. Winter with no metering data;
- 487 2. Spring with one week of metering data transformed in LI;
- 488 3. Summer with four weeks of metering data transformed in LP;
- 489 4. Autumn with eight weeks of metering data transformed in LP.

490 For the classification of the Winter profiles without any smart meter-
491 ing data Table 10 presents the variables selected by the filter FS algorithm
492 (regression analysis) and Figure 14 pictures the rate of selection of the vari-
493 ables selected by the forward FS throughout the cross-validation process.
494 A maximum mean accuracy of 39% is achieved with the features selected
495 by filter FS. With the forward FS the number of features is reduced from

16 to 9 and 4, respectively for LR and DT, achieving a better accuracy for DT (37.4% with forward and 36.3% with filter FS) and slightly worst with LR (37.3%). The variables selected by forward FS with LR modelling are mainly age and employment. `heat_solidfuel`, `tumble_dryer` and `electric_cooker` are also selected in more than half of the cross-validation folds. The variable selected by forward FS with DT modelling is mainly age. `heat_electricity_plugin` and `electric_cooker` are also selected in the more than half of the cross-validation folds. The age, employment, type of heating and the use of electric cooking appliances are the features which can be used as indicators to separate customers with different consumption profiles.

For the classification of the Spring profiles with one week of smart metering data, translated by LI, Table 11 presents the variables selected by the filter FS algorithm and Figure 17 pictures the rate of selection of the variables selected by the forward FS throughout the cross-validation process. A maximum mean accuracy of 56.5% is achieved with the features selected by filter FS. With the forward FS the number of features is reduced from 20 to 9 and 5, respectively for LR and DT, achieving slightly worst accuracies. The variables selected by forward FS with LR modelling are mainly the five LI (i_1, \dots, i_5) and `washing_machine`. The variables selected by forward FS with DT modelling are mainly three LI (i_1, i_3, i_4), indicating that the load factor, night impact and lunch impact are the LI features which can be used as indicators to separate customers with different consumption profiles.

For the classification of the Summer profiles with four weeks of smart metering data, translated by LP, Table 12 presents the variables selected by the filter FS algorithm and Figure 15 pictures the rate of selection of the variables selected by the forward FS throughout the cross-validation process. A maximum mean accuracy of 73.3% is achieved with the features selected by filter FS. With the forward FS the number of features is reduced from 30 to 16 and 5, respectively for LR and DT, achieving slightly worst accuracies (71.7% and 64.9%). The variables selected by forward FS with LR modelling are mainly multiple LP features ($l_1, l_2, l_7, l_{11}, l_{16}, l_{18}, l_{22}, l_{23}, l_{24}$) and `washing_machine`. The variables selected by forward FS with DT modelling are mainly LP features ($l_2, l_{12}, l_{15}, l_{23}$). The consumption behaviour translated by LP features distributed throughout the day in combination with the number of washing machines in the customers household can be used as indicators to separate customers with different consumption profile.

For the classification of the Autumn profiles with eight weeks of smart

metering data, translated by LP, Table 13 presents the variables selected by the filter FS algorithm and Figure 16 pictures the rate of selection of the variables selected by the forward FS throughout the cross-validation process. A maximum mean accuracy of 81.6% is achieved with the features selected by filter FS. With the forward FS the number of features is reduced from 32 to 16 and 8, respectively for LR and DT, achieving worst accuracies (77.9% and 70.4%). The variables selected by forward FS with LR modelling are mainly multiple LP features ($l_8, l_{10}, l_{12}, l_{13}, l_{14}, l_{15}, l_{17}, l_{20}, l_{22}, l_{23}, l_{24}$) and washing_machine. The variables selected by forward FS with DT modelling are mainly LP features ($l_2, l_3, l_5, l_{21}, l_{23}$). The consumption behaviour translated by LP features distributed throughout the day in combination with the number of washing machines in the customers household can be used as indicators to separate customers with different consumption profile.

Notice the LR results having a high standard deviation of the accuracy, such as the results for ten weeks of metering data for Winter and Spring with no FS, using LP metering features. These result due the inappropriate convergence of the optimization method for LR training. Using forward FS this problem is avoided.

Based on the results, the five most important variables or questions an utility should ask customers on sign-up are:

1. What is the customer employment status;
2. How old the customer is;
3. How many dishwashers are used in the clients household;
4. How many electric cookers are used in the clients household;
5. How many washing machines are used in the clients household.

6. Conclusions

The integration of smart metering in the power grid enables a detailed analysis of the consumption behaviour of electricity customers. Knowledge on the typical consumption profiles of customers and the main drivers of consumption are extremely valuable for decision makers in the utility industry and policy. The engagement and education of consumers is seen as a key task in order to successfully reap the potential benefits of the smart grid [41]. The daily routines and the social context of consumers needs to be correctly taken into account to efficiently plan and target the correct groups

568 for potential DSM programs and create incentives for consumers to act with
569 regard towards sustainability.

570 The proposed process is a contribution for enabling the modelling of inter-
571 pretable classifiers to predict the consumption profile group of new customers
572 using smart metering data and survey responses. It enables the discovery of
573 the drivers of consumption profiles, e.g., which characteristics of customers
574 are able to translate consumption behaviour differences. This can contribute
575 to the better engagement of consumers and development of measures to in-
576 crease efficiency in the power grid.

577 An application, based on the data from more than three thousand resi-
578 dential electricity customers from Ireland, shows the viability of the proposed
579 methods. Without any metering data the LR is able to correctly classify up
580 to 39% of the customers which is significantly better than randomly insert-
581 ing the customer in one of the four customer groups (with four customer
582 groups). With the growth of available smart metering data the simulations
583 show an increase in accuracy achieving up to 60%, 70% and 80% accuracy,
584 respectively, with 1, 4 and 8 weeks of data.

585 The forward FS results pictured are easily interpreted and resulted in
586 the discovery of the most important features when grouping electricity cus-
587 tomers by their representative consumption profile. [For the Irish population](#)
588 [studied in the paper, information on the representative consumption profile](#)
589 [throughout all the day results in the highest classification accuracy. A low](#)
590 [number of shape indices is not suitable to accurately classify new electricity](#)
591 [customers.](#) The number of washing machines in the customers households is
592 revealed to be a very important feature in the classification task, seemingly
593 being the most influencing feature to the considerable increase of accuracy
594 from the use of survey features added to the smart metering features.

595 **Acknowledgements**

596 This work was supported by FCT, through IDMEC, under LAETA,
597 project UID/EMS/50022/2013 and SusCity (MITP-TB/CS/0026/2013). The
598 work of J. L. Viegas was supported by the PhD in Industry Scholarship
599 SFRH/BDE/95414/2013 from FCT and Novabase. S. M. Vieira acknowl-
600 edges support by Program Investigador FCT (IF/00833/2014) from FCT,
601 co-funded by the European Social Fund (ESF) through the Operational Pro-
602 gram Human Potential (POPH).

603 **References**

- 604 [1] Exxon Mobil Corporation, The outlook for energy: a view to 2040.
- 605 [2] OECD, ICT applications for the smart grid: opportunities and policy
606 implications, OECD Digital Economy Papers (190).
- 607 [3] D. S. Markovic, D. Zivkovic, I. Branovic, R. Popovic, D. Cvetkovic,
608 Smart power grid and cloud computing, Renewable and Sustainable En-
609 ergy Reviews 24 (2013) 566–577.
- 610 [4] U.S. Energy Information Administration, Annual electric power industry
611 report, Tech. rep. (2013).
612 URL <http://www.eia.gov/electricity/data/eia861/>
- 613 [5] Commission Européenne, Benchmarking smart metering deployment in
614 the EU-27 with a focus on electricity (2014).
- 615 [6] R. Granell, C. J. Axon, D. C. Wallom, Clustering disaggregated load
616 profiles using a Dirichlet process mixture model, Energy Conversion and
617 Management 92 (2015) 507–516.
- 618 [7] P. R. Jota, V. R. Silva, F. G. Jota, Building load management using
619 cluster and statistical analyses, International Journal of Electrical Power
620 & Energy Systems 33 (8) (2011) 1498–1505.
- 621 [8] I. Benítez, A. Quijano, J.-L. Díez, I. Delgado, Dynamic clustering seg-
622 mentation applied to load profiles of energy consumption from Spanish
623 customers, International Journal of Electrical Power & Energy Systems
624 55 (2014) 437–448.
- 625 [9] M. Misiti, Y. Misiti, G. Oppenheim, Optimized clusters for disaggre-
626 gated electricity load forecasting, REVSTAT - Statistical Journal 8 (2)
627 (2010) 105–124.
- 628 [10] F. Andersen, H. Larsen, T. Boomsma, Long-term forecasting of hourly
629 electricity load: Identification of consumption profiles and segmentation
630 of customers, Energy Conversion and Management 68 (2013) 244–252.
- 631 [11] H. R. Sadeghi Keyno, F. Ghaderi, a. Azade, J. Razmi, Forecasting elec-
632 tricity consumption by clustering data in order to decline the periodic

- 633 variable's affects and simplification the pattern, *Energy Conversion and*
634 *Management* 50 (3) (2009) 829–836.
- 635 [12] G. Chicco, I. S. Ilie, Support vector clustering of electrical load pattern
636 data, *IEEE Transactions on Power Systems* 24 (3) (2009) 1619–1628.
- 637 [13] N. Mahmoudi-Kohan, M. P. Moghaddam, M. Sheikh-El-Eslami, An an-
638 nual framework for clustering-based pricing for an electricity retailer,
639 *Electric Power Systems Research* 80 (9) (2010) 1042–1048.
- 640 [14] J. J. López, J. a. Aguado, F. Martín, F. Muñoz, a. Rodríguez, J. E. Ruiz,
641 Hopfield-K-Means clustering algorithm: a proposal for the segmentation
642 of electricity customers, *Electric Power Systems Research* 81 (2) (2011)
643 716–724.
- 644 [15] V. Figueiredo, F. Rodrigues, Z. Vale, J. Gouveia, An electric energy
645 consumer characterization framework based on data mining techniques,
646 *IEEE Transactions on Power Systems* 20 (2) (2005) 596–602.
- 647 [16] T. Räsänen, D. Voukantsis, H. Niska, K. Karatzas, M. Kolehmainen,
648 Data-based method for creating electricity use load profiles using large
649 amount of customer-specific hourly measured electricity use data, *Ap-
650 plied Energy* 87 (11) (2010) 3538–3545.
- 651 [17] L. Hernández, C. Baladrón, J. Aguiar, B. Carro, A. Sánchez-Esguevillas,
652 Classification and clustering of electricity demand patterns in industrial
653 parks, *Energies* 5 (12) (2012) 5215–5228.
- 654 [18] F. Rodrigues, J. Duarte, V. Figueiredo, Z. Vale, M. Cordeiro, A com-
655 parative analysis of clustering algorithms applied to load profiling, in:
656 *Machine Learning and Data Mining in Pattern Recognition*, Springer,
657 2003, pp. 73–85.
- 658 [19] X. Zhang, C. Sun, Dynamic intelligent cleaning model of dirty electric
659 load data, *Energy Conversion and Management* 49 (4) (2008) 564–569.
660 doi:10.1016/j.enconman.2007.08.007.
- 661 [20] A. Mutanen, M. Ruska, Customer classification and load profiling
662 method for distribution systems, *IEEE Transactions on Power Deliv-
663 ery* 26 (3) (2011) 1755–1763.

- 664 [21] T. F. Sanquist, H. Orr, B. Shui, A. C. Bittner, Lifestyle factors in U.S.
665 residential electricity consumption, *Energy Policy* 42 (2012) 354–364.
- 666 [22] A. Kavousian, R. Rajagopal, M. Fischer, Determinants of residential
667 electricity consumption: using smart meter data to examine the effect
668 of climate, building characteristics, appliance stock, and occupants’ be-
669 havior, *Energy* 55 (2013) 184–194.
- 670 [23] M. Bedir, E. Hasselaar, L. Itard, Determinants of electricity consump-
671 tion in Dutch dwellings, *Energy and Buildings* 58 (2013) 194–207.
- 672 [24] B. Sütterlin, T. a. Brunner, M. Siegrist, Who puts the most energy into
673 energy conservation? A segmentation of energy consumers based on
674 energy-related behavioral characteristics, *Energy Policy* 39 (12) (2011)
675 8137–8152.
- 676 [25] F. Fusco, M. Wurst, J. W. Yoon, Mining residential household informa-
677 tion from low-resolution smart meter data, in: 21st International Con-
678 ference on Pattern Recognition (ICPR), IEEE, 2012, pp. 3545–3548.
- 679 [26] C. Beckel, L. Sadamori, T. Staake, S. Santini, Revealing household char-
680 acteristics from smart meter data, *Energy* 78 (2014) 397–410.
- 681 [27] T. K. Wijaya, T. Ganu, D. Chakraborty, K. Aberer, D. P. Seetharam,
682 Consumer segmentation and knowledge extraction from smart meter
683 and survey data, in: SIAM International Conference on Data Mining
684 (SDM14), 2014.
- 685 [28] J. D. Rhodes, W. J. Cole, C. R. Upshaw, T. F. Edgar, M. E. Webber,
686 Clustering analysis of residential electricity demand profiles, *Applied*
687 *Energy* 135 (2014) 461–471. doi:10.1016/j.apenergy.2014.08.111.
- 688 [29] S. Ramos, J. M. Duarte, F. J. Duarte, Z. Vale, A data-mining-based
689 methodology to support MV electricity customers characterization, *En-
690 ergy and Buildings* 91 (2015) 16–25.
- 691 [30] R. N. Bolton, A Dynamic Model of the Duration of the Customer’s Re-
692 lationship with a Continuous Service Provider: The Role of Satisfaction,
693 *Marketing Science* 17 (1) (1998) 45–65. doi:10.1287/mksc.17.1.45.

- 694 [31] P. C. Verhoef, Understanding the effect of customer relation-
 695 ship management efforts on customer retention and customer
 696 share development, *Journal of Marketing* 67 (4) (2003) 30–45.
 697 doi:10.1509/jmkg.67.4.30.18685.
- 698 [32] ISSDA, Data from the Commission for Energy Regulation -
 699 www.ucd.ie/issda.
- 700 [33] M. Welsch, M. Howells, M. Bazilian, J. DeCarolis, S. Hermann,
 701 H. Rogner, Modelling elements of smart grids: enhancing the OSe-
 702 MOSYS (open source energy modelling system) code, *Energy* 46 (1)
 703 (2012) 337–350.
- 704 [34] U.S. Department of Energy, The smart grid: an introduction, Tech. rep.
 705 (2008).
- 706 [35] International Energy Agency, Technology roadmap: smart grids, Tech.
 707 rep. (2011).
- 708 [36] ETP SmartGrids, European technology platform smart grids: vision
 709 and strategy for Europe’s electricity networks of the future, Tech. rep.
 710 (2006).
- 711 [37] A. Battaglini, J. Lilliestam, A. Haas, A. Patt, Development of supers-
 712 mart grids for a more efficient utilisation of electricity from renewable
 713 sources, *Journal of Cleaner Production* 17 (10) (2009) 911–918.
- 714 [38] C. Felix, M. Ardelean, J. Vasiljevska, A. Mengolini, G. Fulli,
 715 E. Amoiralis, M. S. Jiménez, C. Filiou, Smart grid projects outlook
 716 2014, European Commision, JRC Science and Policy Reports.
- 717 [39] A. Faruqui, D. Harris, R. Hledik, Unlocking the €53 billion savings
 718 from smart meters in the eu: How increasing the adoption of dynamic
 719 tariffs could make or break the eu’s smart grid investment, *Energy Policy*
 720 38 (10) (2010) 6222–6231.
- 721 [40] A. J. Conejo, J. M. Morales, L. Baringo, Real-time demand response
 722 model, *Smart Grid, IEEE Transactions on* 1 (3) (2010) 236–242.
- 723 [41] V. Giordano, F. Gangale, G. Fulli, M. Sánchez, J. Dg, I. Onyeji,
 724 A. Colta, I. Papaioannou, A. Mengolini, C. Alecu, T. Ojala, I. Maschio,

- 725 Smart grid projects in Europe : lessons learned and current develop-
726 ments, European Commision: JRC Scientific and Policy Reports.
- 727 [42] Institute of Communication & Computer Systems of the National Tech-
728 nical University of Athen ICCS-NTUA for the European Commission,
729 Study on cost benefit analysis of Smart Metering Systems in EU Member
730 States - Final Report.
- 731 [43] J. L. Viegas, S. M. Vieira, R. Melício, V. M. F. Mendes, J. a. M. C.
732 Sousa, Electricity demand profile prediction based on household char-
733 acteristics, in: Proceedings of the 12th International Conference on the
734 European Energy Market, 2015.
- 735 [44] J. MacQueen, Some methods for classification and analysis of multi-
736 variate observations, in: Proceedings of the fifth Berkeley symposium
737 on mathematical statistics and probability, Vol. 1, Oakland, CA, USA.,
738 1967, pp. 281–297.
- 739 [45] G. Chicco, Overview and performance assessment of the clustering meth-
740 ods for electrical load pattern grouping, *Energy* 42 (1) (2012) 68–80.
741 doi:10.1016/j.energy.2011.12.031.
742 URL <http://dx.doi.org/10.1016/j.energy.2011.12.031>
- 743 [46] T. Warren Liao, Clustering of time series data - A survey, *Pattern Recog-*
744 *nition* 38 (2005) 1857–1874.
- 745 [47] O. Arbelaiz, I. Gurrutxaga, J. Muguerza, J. M. Pérez, I. n. Perona, An
746 extensive comparative study of cluster validity indices, *Pattern Recog-*
747 *nition* 46 (1) (2013) 243–256.
- 748 [48] C. Dunn, A fuzzy relative of the ISODATA process and its use in detect-
749 ing compact well-separated clusters, *Journal of Cybernetics* 3 (3) (1973)
750 32–57.
- 751 [49] D. L. Davies, D. W. Bouldin, A cluster separation measure, *IEEE Trans-*
752 *actions on Pattern Analysis and Machine Intelligence* (2) (1979) 224–
753 227.
- 754 [50] P. Rousseeuw, Silhouettes: A graphical aid to the interpretation and
755 validation of cluster analysis, *Journal of Computational and Applied*
756 *Mathematics* 20 (1987) 53–65.

- 757 [51] T. Hastie, R. Tibshirani, J. Friedman, The elements of statistical learn-
758 ing: data mining, inference, and prediction, Springer.[Online book],
759 2008.
- 760 [52] D. C. Montgomery, E. A. Peck, G. Vining, Introduction to Linear Re-
761 gression Analysis, 5th Edition, Wiley, 2012.
- 762 [53] M. R. Berthold, C. Borgelt, F. Höppner, F. Klawonn, Guide to Intel-
763 ligent Data Analysis: How to Intelligently Make Sense of Real Data,
764 Springer, 2010.
- 765 [54] S. M. Vieira, J. a. M. Sousa, T. a. Runkler, Two cooperative ant colonies
766 for feature selection using fuzzy models, Expert Systems with Applica-
767 tions 37 (4) (2010) 2714–2723. doi:10.1016/j.eswa.2009.08.026.
- 768 [55] Y. Saeys, I. Inza, P. Larrañaga, A review of feature selection techniques
769 in bioinformatics, Bioinformatics 23 (19) (2007) 2507–2517.
- 770 [56] J. Kittler, Feature set search algorithms, Pattern recognition and signal
771 processing (1978) 41–60.
- 772 [57] G. Chicco, R. Napoli, P. Postolache, M. Scutariu, C. Toader, Customer
773 characterization options for improving the tariff offer, IEEE Transac-
774 tions on Power Systems 18 (1) (2003) 381–387.
- 775 [58] J. A. Hanley, B. J. McNeil, The meaning and use of the area under a
776 receiver operating characteristic (ROC) curve., Radiology 143 (4) (1982)
777 29–36. doi:10.1148/radiology.143.1.7063747.

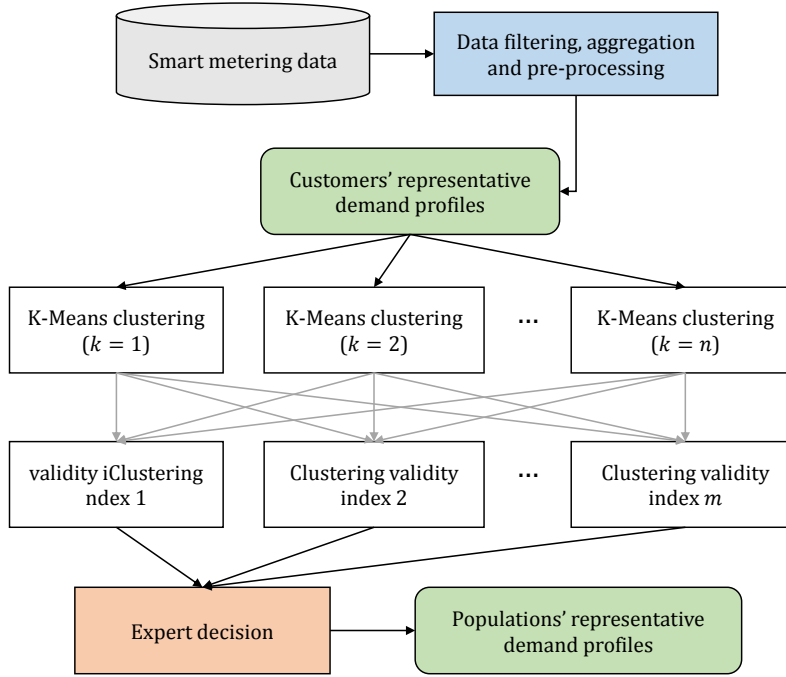
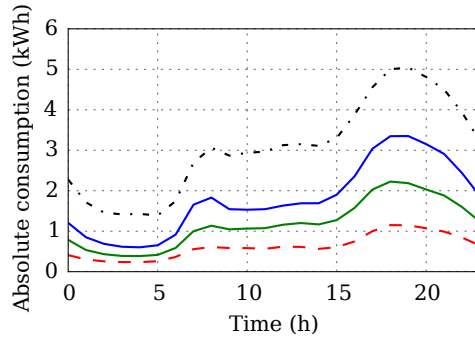
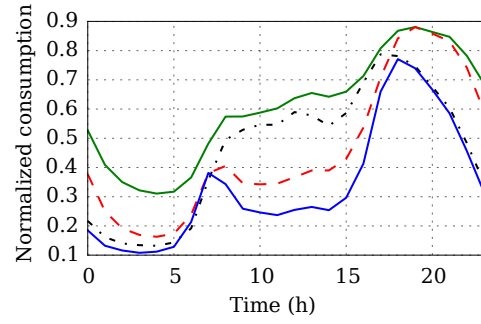


Figure 1: Generation of populations representative consumption profiles.



(a) Absolute consumption



(b) Customer normalized consumption

Figure 2: Example of populations representative consumption profiles using absolute and customer normalized consumption (resulting cluster centroids).

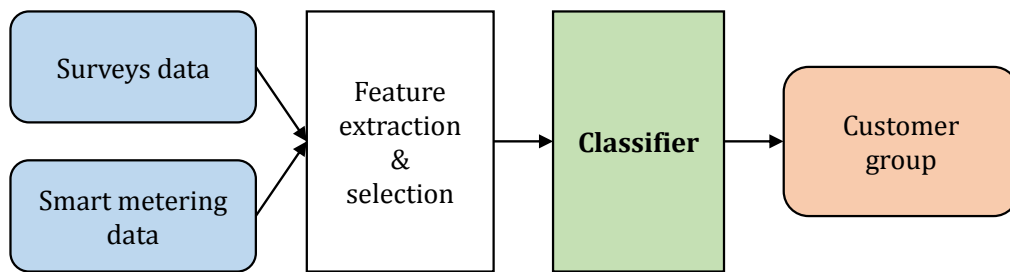


Figure 3: Electricity customer classifier.

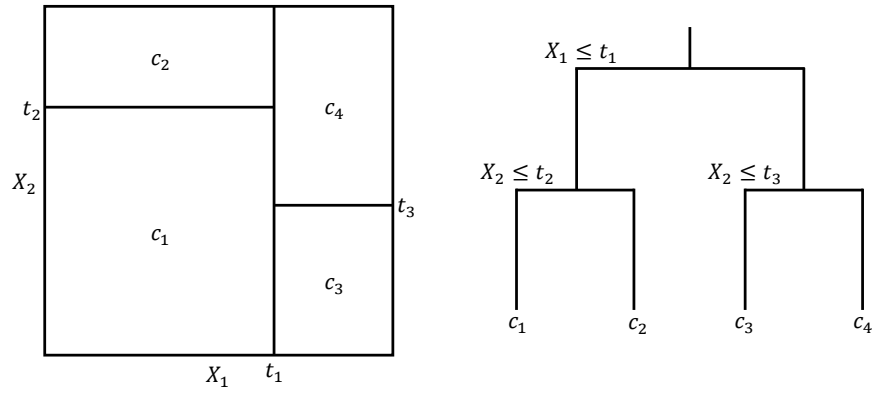


Figure 4: Left: data partitioned in four categories by binary splitting. Right: CART tree corresponding to the partition.

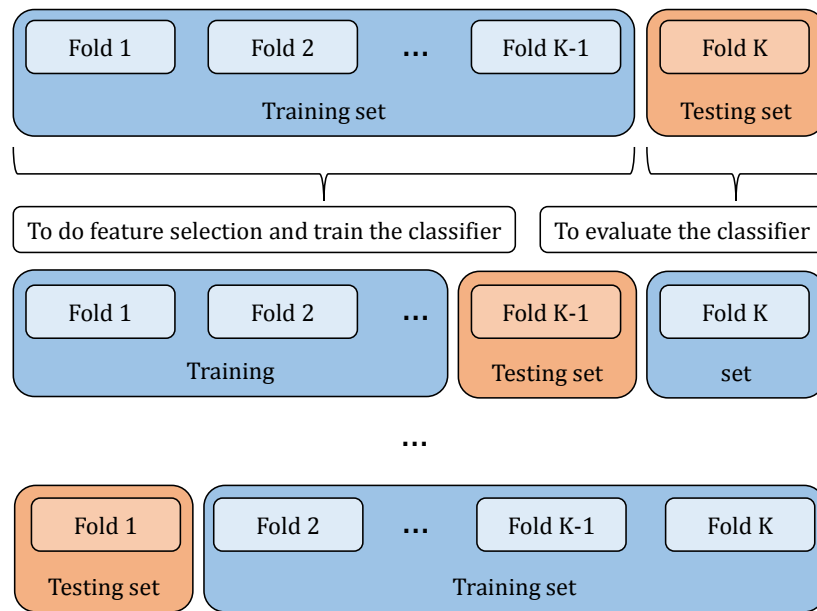


Figure 5: K-fold cross-validation.

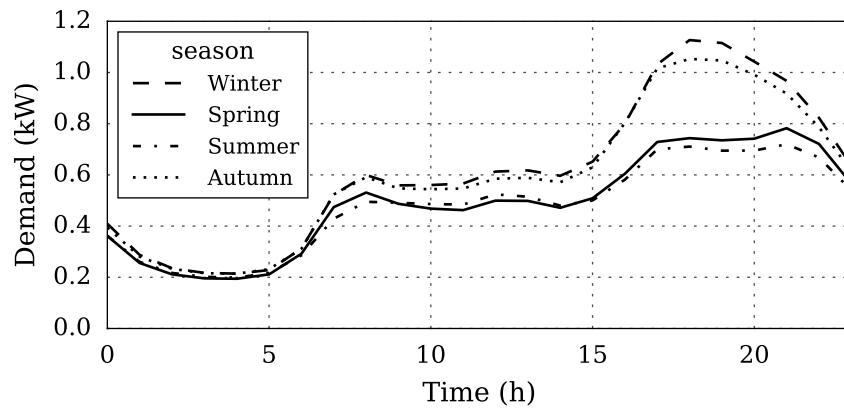


Figure 6: Hourly aggregated mean seasonal consumption of all customers.

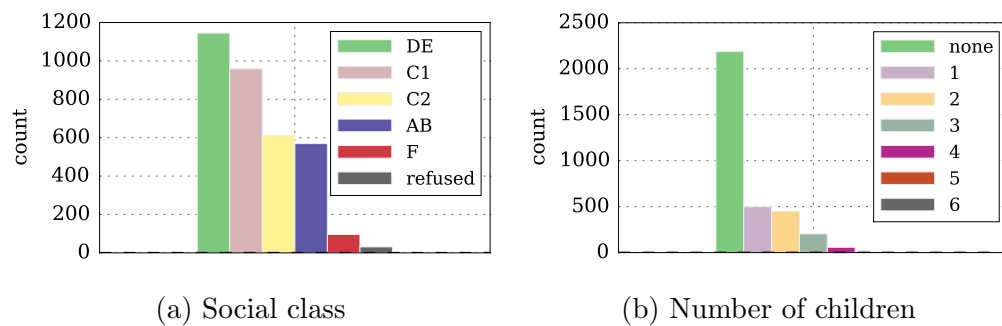


Figure 7: Distribution of the households for two survey responses.

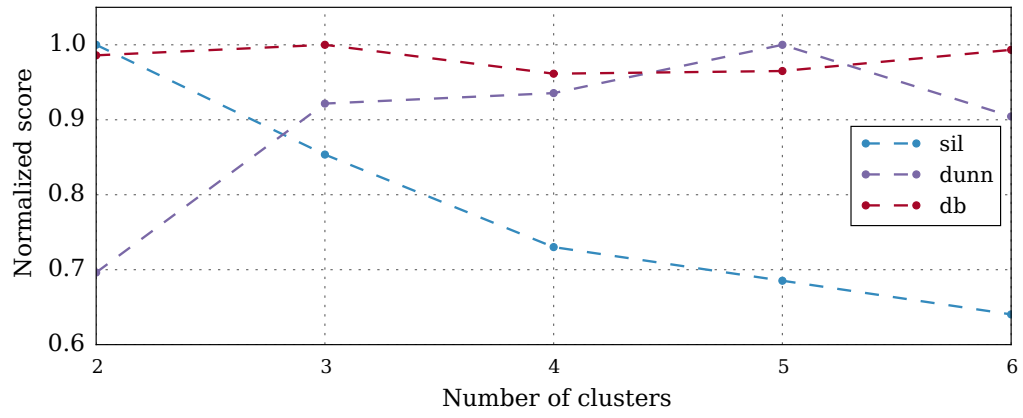


Figure 8: CVI for different number of clusters for the Winter consumption profiles.

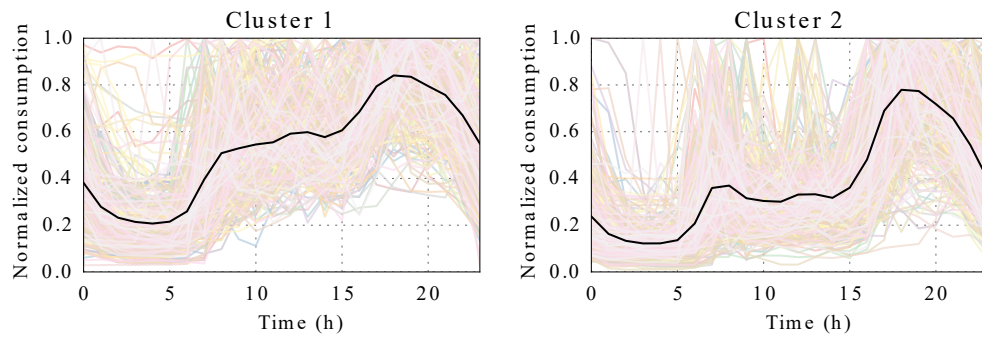


Figure 9: Winter clustering results with two clusters.

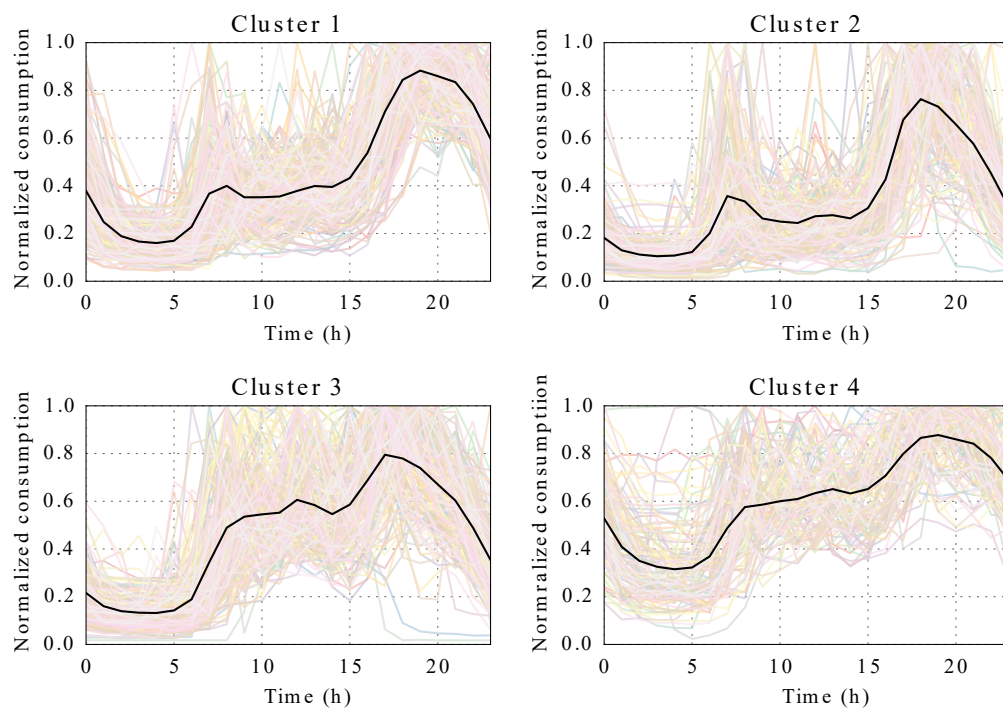


Figure 10: Winter clustering results with four clusters.

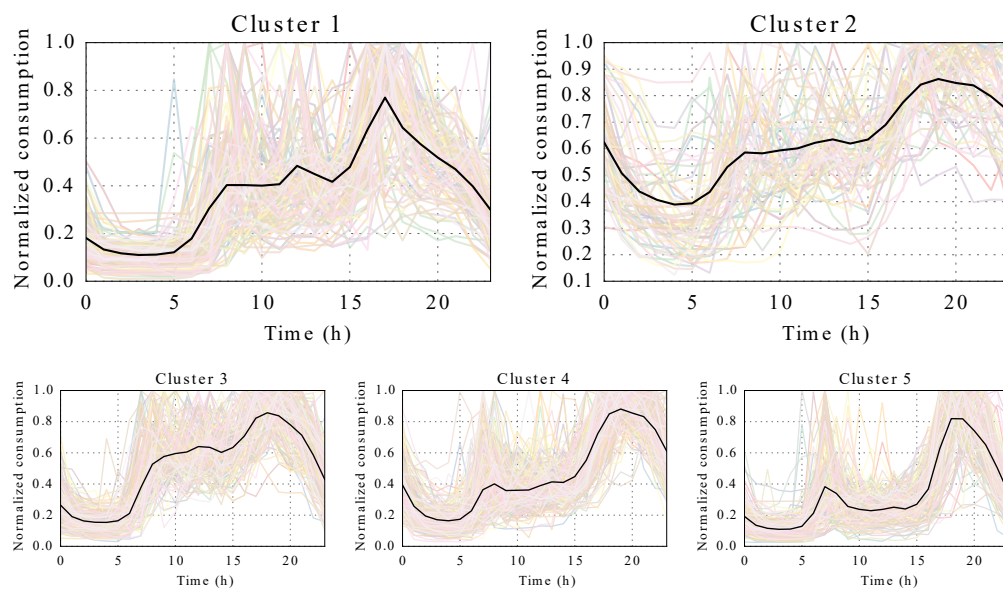
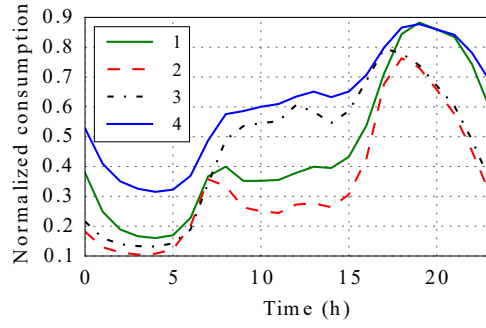
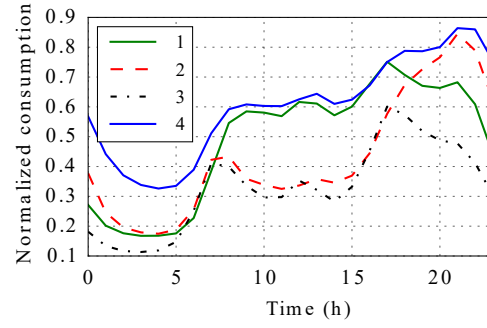


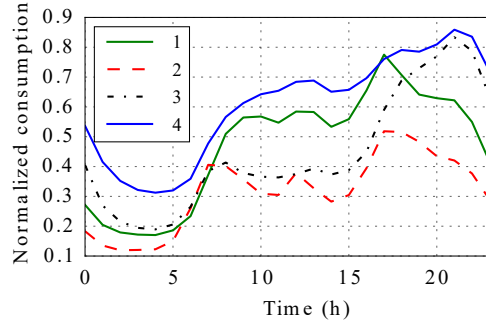
Figure 11: Winter clustering results with five clusters.



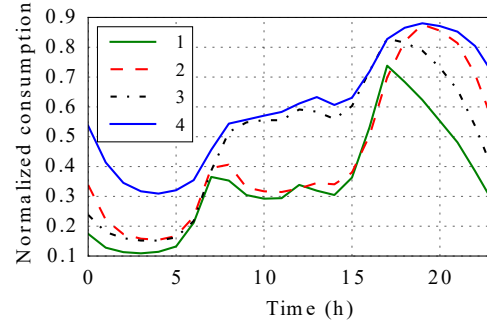
(a) Winter



(b) Spring

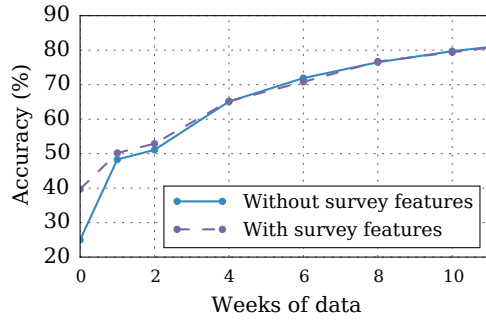


(c) Summer

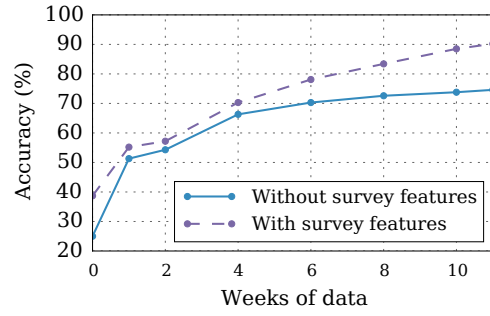


(d) Autumn

Figure 12: Populations representative consumption profiles.



(a) Metering data features: LI



(b) Metering data features: LP

Figure 13: LR classifier accuracy using filter FS with and without the survey features for Winter profiles.

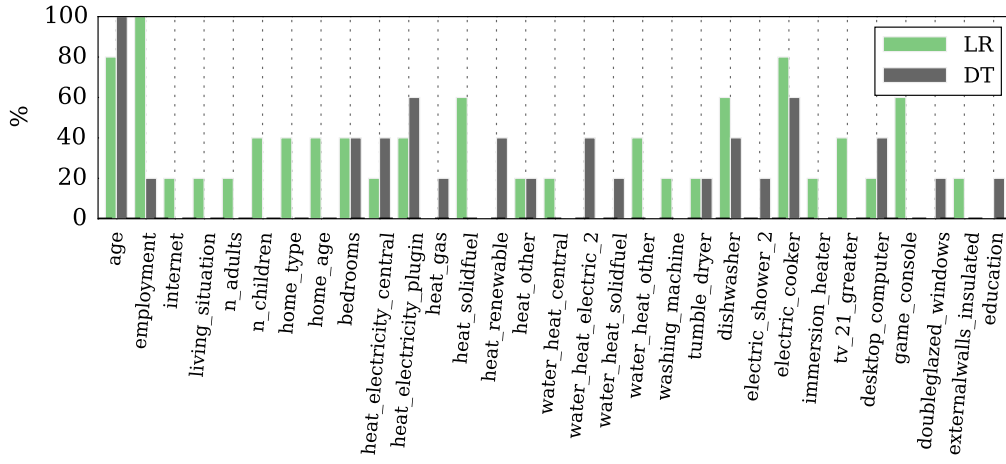


Figure 14: Forward FS for Winter with no metering data: rate of selection of features throughout the cross-validation process.

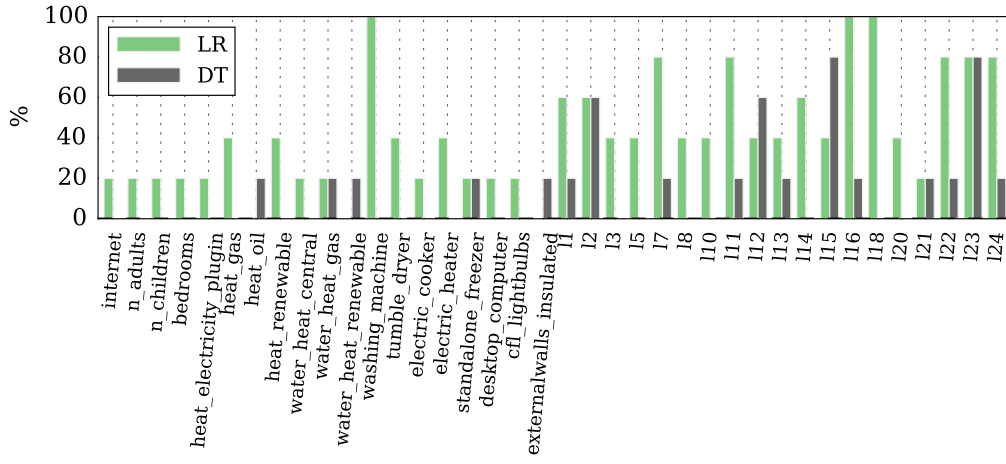


Figure 15: Forward FS for Summer with 4 weeks metering data (LP): rate of selection of features throughout the cross-validation process.

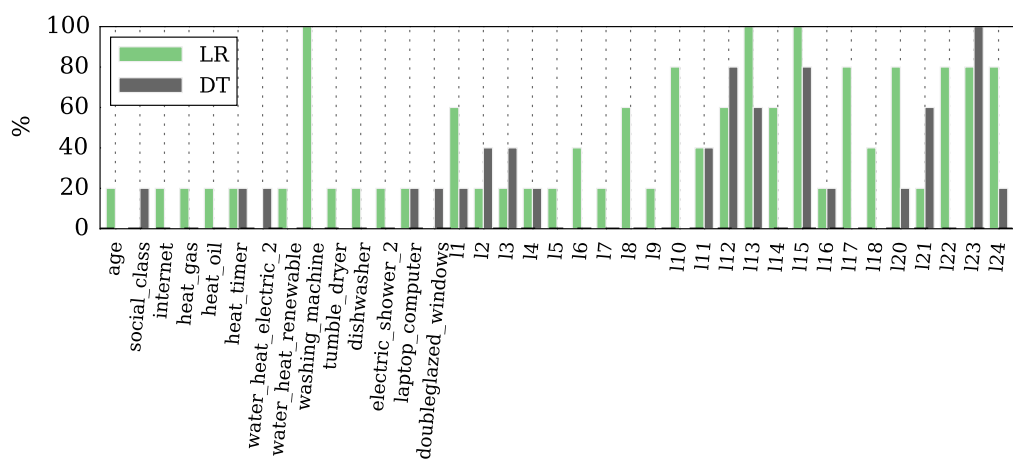


Figure 16: Forward FS for Autumn with 8 weeks metering data (LP): rate of selection of features throughout the cross-validation process.

Nomenclature			
<i>Acronyms</i>			
AMI	Advanced metering infrastructure	r	number of smart metering data features
EU	European Union	t	number of survey features
DSM	Demand side management	X	feature dataset of all customers
CVI	Clustering validity index	μ_i	i th consumption profile of the population
LP	Load profile		
LI	Load indexes	\mathbf{S}	set of the groups of customers
FS	Feature selection	S_i	i th clustered group of customers
LR	Logistic regression		
DT	Decision tree	J	number of clusters/customer groups
CER	Commission for Energy Regulation	$d_e(\mathbf{v}_1, \mathbf{v}_1)$	euclidean distance
ISSDA	Irish Social Science Data Archive	$D(\mathbf{S})$	Dunn index
<i>Symbols</i>			
\mathbf{x}_i	feature vector of customer i	$DB(S)$	Davis Bouldin index
\mathbf{x}_i^m	customer i smart metering data features	$Sil(\mathbf{S})$	Silhouette index
\mathbf{x}_i^s	customer i surveys features	y	categorical variable representing a group
N	number of customers	i_1, i_2, \dots, i_5	load indices
p	dimension of feature vector	$P_{max/min/av}$	maximum, minimum and average consumption
		l_1, l_2, \dots, l_{24}	load profile

Table 1: Survey features I: respondent

Feature	Description: {responses}
sex	Sex of respondent: {male, female}
age	Age of respondent in years: {18-25, 26-35, 36-45, 46-55, 56-65, 65 or more, refused}
employment	employment status of respondent: {Employee, self-employed, unemployed}
social_class	Social class of respondent: {AB, C1, C2, DE, F, refused}
education	Education level of respondent: {none, primary, secondary to intermediate cert junior cert level, secondary to leaving cert level, third level, refused}
income	Income of respondent before tax in euro: {0-15k, 15k-30k, 30k-50k, 50k-75k, 75k or more, refused}

Table 2: Survey features II: household

Feature	Description: {responses}
home_type	Household type: {apartment, semi-detached, detached, terraced, bungalow}
home_age	Household age in years: {0-4, 5-9, 10-29, 30-74, 75 or more}
bedrooms	Number of bedrooms : {1, 2, 3, 4, 5 or more, refused}
clf_lighbulbs	Fraction of CLF light bulbs: {none, about a quarter, about half, about three quarters}
doulegazed_windows	Fraction of doubleglazed windows: {none, about a quarter, about half, about three quarters}
attic_insulated	Presence and age of attic insulation: {yes (last 5 years), yes, no, don't know}
externalwalls_insuled	Presence and age of insulation of external walls: {yes, no, don't know}
internet	Internet connection in the household: {yes, no}

Table 3: Survey features III: heating

Feature	Description: {responses}
heat_electricity_central	Central electric heating : {yes, no}
heat_gas	Gas heating : {yes, no}
heat_oil	Oil heating : {yes, no}
heat_solidfuel	Solid fuel heating : {yes, no}
heat_renewable	Renewable energy heating : {yes, no}
heat_other	Other type of heating : {yes, no}
heat_timer	Use of heating timer : {yes, no}
water_heat_central	Central water heating : {yes, no}
water_heat_electric	Electric water heating: {yes, no}
water_heat_gas	Gas water heating: {yes, no}
water_heat_oil	Oil water heating: {yes, no}
water_heat_solidfuel	Solid fuel water heating: {yes, no}
water_heat_renewable	Renewable water heating: {yes, no}
water_heat_other	Other water heating source : {yes, no}

Table 4: Survey features IV: appliances

Feature	Description: {responses}
washing_machine	Number of washing machines : {0, 1, 2, 3 or more}
tumble_dryer	Number of tumble dryers : {0, 1, 2, 3 or more}
dishwasher	Number of dishwashers : {0, 1, 2, 3 or more}
electric_shower	Number of electric showers : {0, 1, 2, 3 or more}
electric_cooker	Number of electric cookers : {0, 1, 2, 3 or more}
electric_heater	Number of electric heaters : {0, 1, 2, 3 or more}
standalone_freezer	Number of standalone freezers : {0, 1, 2, 3 or more}
water_pump	Number of water pumps : {0, 1, 2, 3 or more}
immersion_heater	Number of immersion heaters : {0, 1, 2, 3 or more}
tv_21_less	Numbers of TVs with 21 or less inches: {0, 1, 2, 3, 4 or more}
tv_21_greater	Number of TVs with more than 21 inches: {0, 1, 2, 3, 4 or more}
desktop_computer	Number of desktop computers: {0, 1, 2, 3, 4 or more}
laptop_computer	Number of laptop computers: {0, 1, 2, 3, 4 or more}
game_console	Number of game consoles: {0, 1, 2, 3, 4 or more}

Table 5: Distribution of customers between the different clusters for the four seasons

Cluster	Winter	Spring	Summer	Autumn
1	30.93%	26.25%	26.14%	20.34%
2	25.50%	31.89%	18.83%	31.47%
3	28.17%	21.53%	27.17%	29.19%
4	15.39%	20.33%	27.86%	18.99%

Table 6: Normalized indices to characterize electricity customers' behaviour

Parameter	Definition	Periods
Daily P_{av}/P_{max}	$i_1 = P_{av,day}/P_{max,day}$	1 day
Daily $P_{min,day}/P_{max,day}$	$i_2 = P_{min,day}/P_{max,day}$	1 day
Night impact	$i_3 = 1/3P_{av,night}/P_{av,day}$	1 day and 8 h night (from 23h to 06h)
Lunch impact	$i_4 = 1/8P_{av,lunch}/P_{av,day}$	1 day and 3 h lunch from (12h to 15h)
Daily P_{min}/P_{av}	$i_5 = P_{min,day}/P_{av,day}$	1 day

Table 7: Smart metering data features used for classification

Smart metering data features		
Load indices (LI)	Normalized indices to characterize electricity costumers' behaviour.	i_1, i_2, i_3, i_4, i_5
Load profile (LP)	Normalized mean hourly aggregated consumption.	l_1, l_2, \dots, l_{24}

Table 8: Mean 10-fold cross-validation accuracy of classifiers using load indices as metering data features (number of selected features)

Smart metering data features: Load indices					
W	Model	Winter	Spring	Summer	Autumn
No FS					
0	LR	39.2±0.8 (47)	37.6±1.1 (47)	37.0±1.9 (47)	38.5±0.7 (47)
	DT	36.0±1.5 (47)	34.7±1.1 (47)	33.8±2.4 (47)	36.7±1.7 (47)
1	LR	45.3±6.9 (52)	54.9±1.1 (52)	53.1±1.5 (52)	53.4±1.3 (52)
	DT	46.5±1.7 (52)	53.3±1.6 (52)	51.8±1.5 (52)	51.4±2.0 (52)
4	LR	64.9±1.8 (52)	64.6±1.3 (52)	65.7±2.1 (52)	62.0±1.4 (52)
	DT	62.8±0.8 (52)	62.4±2.7 (52)	63.3±2.7 (52)	59.3±1.9 (52)
8	LR	75.8±1.3 (52)	71.8±0.8 (52)	71.0±0.9 (52)	57.1±19.0 (52)
	DT	73.3±1.3 (52)	70.4±0.5 (52)	69.3±2.3 (52)	67.7±1.6 (52)
10	LR	78.3±1.4 (52)	75.4±0.8 (52)	64.9±19.1 (52)	73.4±1.8 (52)
	DT	75.1±1.0 (52)	72.9±0.8 (52)	71.7±1.3 (52)	72.1±1.8 (52)
Filter FS					
0	LR	38.6±1.8 (17)	36.2±2.3 (18)	35.9±1.1 (17)	34.6±7.6 (23)
	DT	36.7±1.7 (17)	35.7±0.9 (18)	34.1±2.2 (17)	35.1±0.5 (23)
1	LR	49.8±0.8 (21)	56.5±1.7 (20)	53.9±2.3 (21)	53.4±1.3 (19)
	DT	46.3±2.7 (21)	52.8±1.4 (20)	51.0±0.3 (21)	50.6±1.4 (19)
4	LR	58.4±12.4 (26)	65.9±1.3 (18)	66.8±0.3 (16)	62.6±1.7 (19)
	DT	62.1±2.7 (26)	62.0±0.8 (18)	64.2±1.7 (16)	60.1±1.0 (19)
8	LR	76.5±2.4 (19)	72.7±1.8 (17)	72.0±0.5 (17)	59.4±19.7 (26)
	DT	73.8±0.9 (19)	69.5±2.0 (17)	69.8±1.7 (17)	67.7±1.2 (26)
10	LR	79.1±1.5 (17)	76.0±2.0 (19)	75.2±0.7 (15)	74.3±1.6 (22)
	DT	75.3±1.9 (17)	72.6±1.1 (19)	71.9±1.8 (15)	72.1±1.2 (22)
Forward FS					
0	LR	38.2±1.1 (9)	36.7±1.3 (5)	34.8±1.5 (10)	37.7±4.3 (5)
	DT	36.6±1.1 (6)	35.8±0.5 (4)	32.9±1.1 (5)	37.4±4.3 (2)
1	LR	49.5±1.1 (11)	56.0±2.4 (9)	54.3±1.7 (10)	53.0±3.6 (10)
	DT	46.9±1.5 (6)	52.6±0.9 (5)	52.3±0.9 (5)	50.6±2.2 (4)
4	LR	50.7±16.7 (6)	65.5±1.3 (11)	66.3±1.5 (7)	62.5±1.4 (7)
	DT	61.7±1.9 (4)	62.9±2.1 (4)	63.1±1.2 (5)	60.4±0.8 (4)
8	LR	76.4±1.3 (8)	72.1±2.5 (8)	72.2±0.9 (9)	70.7±1.0 (8)
	DT	71.5±0.8 (4)	69.8±1.3 (4)	69.5±1.2 (4)	67.4±1.9 (4)
10	LR	79.2±1.8 (9)	75.6±1.5 (9)	76.0±1.4 (9)	74.3±1.5 (8)
	DT	75.8±1.3 (4)	72.7±2.6 (4)	72.3±1.2 (4)	71.8±0.8 (3)

Table 9: Mean 10-fold cross-validation accuracy of classifiers using the load profile as metering data features (mean number of selected features)

Smart metering data features: Load profile					
W	Model	Winter	Spring	Summer	Autumn
No FS					
0	LR	38.7±2.1 (47)	37.2±2.6 (47)	36.6±2.3 (47)	38.6±0.9 (47)
	DT	36.0±1.3 (47)	34.7±1.2 (47)	34.4±1.8 (47)	35.1±0.8 (47)
1	LR	53.9±0.9 (71)	60.8±2.4 (71)	58.7±2.3 (71)	60.6±1.3 (71)
	DT	48.0±2.6 (71)	54.8±2.2 (71)	52.0±0.9 (71)	53.3±2.0 (71)
4	LR	70.8±1.9 (71)	72.5±1.6 (71)	72.3±1.0 (71)	70.6±1.2 (71)
	DT	63.6±2.5 (71)	65.2±1.7 (71)	65.6±2.1 (71)	64.5±1.7 (71)
8	LR	83.4±1.5 (71)	80.8±1.3 (71)	79.4±1.1 (71)	78.4±1.2 (71)
	DT	74.1±1.6 (71)	72.7±1.5 (71)	70.6±1.8 (71)	71.7±1.9 (71)
10	LR	76.3±22.2 (71)	73.1±24.7 (71)	82.9±1.0 (71)	83.2±1.2 (71)
	DT	76.6±1.2 (71)	74.9±1.0 (71)	73.2±1.1 (71)	75.4±1.8 (71)
Filter FS					
0	LR	39.0±1.2 (16)	37.4±0.9 (17)	35.3±1.3 (18)	38.9±1.4 (16)
	DT	36.3±0.7 (16)	35.4±0.7 (17)	34.2±1.2 (18)	36.5±1.7 (16)
1	LR	53.9±0.8 (29)	60.8±0.6 (28)	59.7±1.1 (32)	60.9±1.5 (28)
	DT	49.0±0.9 (29)	55.0±1.5 (28)	52.6±1.2 (32)	53.6±2.6 (28)
4	LR	62.2±16.5 (40)	72.9±1.6 (29)	73.3±1.0 (30)	71.4±2.2 (29)
	DT	63.9±0.4 (40)	64.1±1.7 (29)	64.9±1.5 (30)	64.3±1.7 (29)
8	LR	83.1±2.4 (32)	81.6±0.8 (32)	79.6±1.8 (34)	78.9±1.9 (33)
	DT	73.4±2.3 (32)	72.8±0.9 (32)	70.8±1.0 (34)	71.4±1.9 (33)
10	LR	88.3±0.9 (37)	86.1±0.6 (42)	83.1±0.9 (37)	83.8±0.8 (38)
	DT	76.3±1.1 (37)	76.5±1.5 (42)	72.4±0.7 (37)	76.8±1.0 (38)
Forward FS					
0	LR	37.3±5.2 (9)	36.8±0.9 (8)	34.1±1.8 (8)	37.8±1.0 (7)
	DT	37.4±2.3 (4)	36.5±0.9 (3)	32.2±1.2 (4)	36.2±2.0 (3)
1	LR	50.9±1.4 (11)	59.1±1.6 (13)	56.4±2.4 (13)	57.9±1.6 (12)
	DT	48.2±1.9 (5)	52.6±1.9 (6)	51.9±2.4 (6)	51.6±1.3 (6)
4	LR	69.8±1.6 (16)	70.3±1.3 (11)	71.7±1.7 (16)	70.2±1.4 (13)
	DT	63.4±1.4 (7)	64.3±1.5 (6)	64.9±1.9 (5)	62.9±2.3 (6)
8	LR	83.4±1.1 (14)	80.8±1.1 (16)	77.9±2.1 (15)	77.8±1.2 (14)
	DT	73.2±1.4 (5)	71.9±2.1 (8)	70.4±2.5 (7)	70.8±1.9 (6)
10	LR	87.3±0.8 (16)	85.2±0.9 (17)	81.1±2.0 (16)	82.8±1.3 (15)
	DT	76.1±1.1 (6)	75.1±2.2 (6)	73.4±2.2 (5)	74.2±1.1 (7)

Table 10: Filter FS for Winter with no metering data: variables found to be significant for at least one of the classifiers of the MNLogit

Filter FS: Winter with no metering data		
age	employment	social_class
living_situation	n_children	bedrooms
water_heat_oil	dishwasher	electric_shower_1
electric_shower_2	electric_cooker	electric_heater
tv_21_greater	desktop_computer	game_console
cfl_lightbulbs	cfl_lightbulbs	cfl_lightbulbs

Table 11: Filter FS for Spring with 1 week metering data (LI): variables found to be significant for at least one of the classifiers of the MNLogit

Filter FS: Spring with one week metering data (LI)		
age	employment	living_situation
n_children	home_type	home_age
bedrooms	heat_solidfuel	water_heat_solidfuel
washing_machine	tumble_dryer	dishwasher
electric_shower_2	electric_cooker	tv_21_less
externalwalls_insulated	education	income
i1	i2	i3
i4	i4	i4

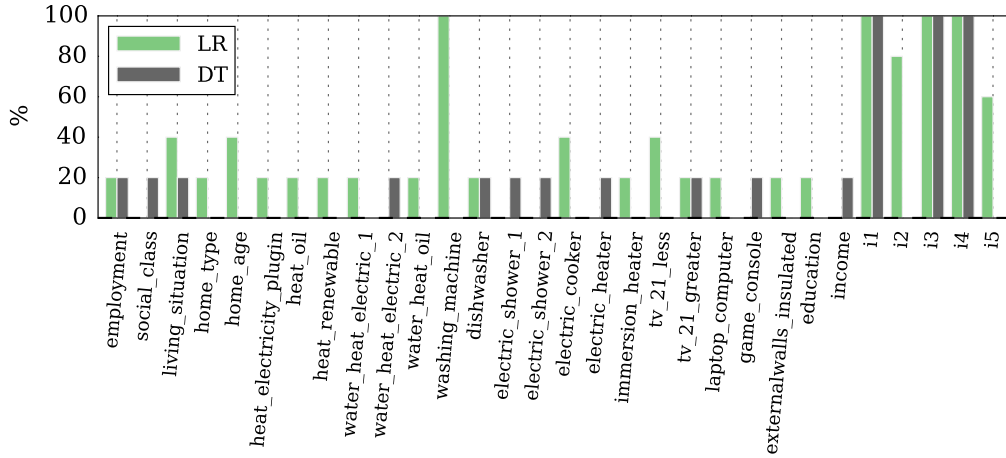


Figure 17: Forward FS for Spring with 1 week metering data (LI): rate of selection of features throughout the cross-validation process.

Table 12: Filter FS for Summer with four weeks metering data (LP): variables found to be significant for at least one of the classifiers of the MNLogit

Filter FS: Summer with four weeks metering data (LP)		
age	social_class	internet
living_situation	n_children	home_type
water_heat_electric_2	water_heat_oil	washing_machine
electric_cooker	standalone_freezer	l1
l3	l8	l9
l10	l11	l12
l13	l14	l15
l16	l17	l18
l19	l20	l21
l22	l23	l24

Table 13: Filter FS for Autumn with eight weeks metering data (LP): variables found to be significant for at least one of the classifiers of the MNLogit

Filter FS: Autumn with eight weeks metering data (LP)		
internet	living_situation	heat_timer
water_heat_electric_2	water_heat_gas	water_heat_oil
washing_machine	tumble_dryer	electric_cooker
game_console	cfl_lightbulbs	attic_insulated
externalwalls_insulated	education	l2
l5	l6	l9
l10	l11	l12
l13	l14	l15
l16	l17	l18
l19	l20	l21
l22	l23	l24