

VTT Technical Research Centre of Finland

Hierarchical Multiplicative Model for Characterizing Residential Electricity Consumption

Kuusela, Pirkko; Norros, Ilkka; Reittu, Hannu; Piira, Kalevi

Published in:
Journal of Energy Engineering

DOI:
[10.1061/\(ASCE\)EY.1943-7897.0000532](https://doi.org/10.1061/(ASCE)EY.1943-7897.0000532)

Published: 01/03/2018

Document Version
Early version, also known as pre-print

[Link to publication](#)

Please cite the original version:
Kuusela, P., Norros, I., Reittu, H., & Piira, K. (2018). Hierarchical Multiplicative Model for Characterizing Residential Electricity Consumption. *Journal of Energy Engineering*, 144(3).
[https://doi.org/10.1061/\(ASCE\)EY.1943-7897.0000532](https://doi.org/10.1061/(ASCE)EY.1943-7897.0000532)



VTT
<http://www.vtt.fi>
P.O. box 1000FI-02044 VTT
Finland

By using VTT's Research Information Portal you are bound by the following Terms & Conditions.

I have read and I understand the following statement:

This document is protected by copyright and other intellectual property rights, and duplication or sale of all or part of any of this document is not permitted, except duplication for research use or educational purposes in electronic or print form. You must obtain permission for any other use. Electronic or print copies may not be offered for sale.

Hierarchical multiplicative model for characterizing residential electricity consumption

Pirkko Kuusela¹ Ph.D. Ilkka Norros² Prof. Hannu Reittu³ Ph.D. and Kalevi Piira⁴

¹VTT Technical Research Centre of Finland Ltd., P.O. Box 1000,
FI-02044 VTT, Finland. pirkko.kuusela@vtt.fi (corresponding author),

²VTT Technical Research Centre of Finland Ltd., P.O. Box 1000,
FI-02044 VTT, Finland. ilkka.norros@vtt.fi,

³VTT Technical Research Centre of Finland Ltd., P.O. Box 1000,
FI-02044 VTT, Finland. hannu.reittu@vtt.fi,

⁴VTT Technical Research Centre of Finland Ltd., P.O. Box 1000,
FI-02044 VTT, Finland. kalevi.piira@vtt.fi

Abstract

This work presents a hierarchical multiplicative framework for modeling the energy consumption of households. The constituents of the model are a lognormally distributed annual consumption, an annual consumption profile at week resolution, a mean weekly consumption profile, and a multiplicative lognormally distributed random variation. Further, the annual and weekly profiles of households are shown to fall naturally into a small number of rather homogeneous groups, identified by the Regular Decomposition method. The framework is adapted to monitor and compare populations of electricity consumers. On the other hand, it provides a convenient way to produce synthetic traces of household energy consumption with similar stochastic properties as measured traces. It is also shown how additional household information can be utilized to predict both the annual consumption and the random variation of the consumption of a household.

Keywords: household electricity consumption, mathematical modeling, clustering, profiles, monitoring

18 INTRODUCTION

19 Local (district level and building embedded) renewable energy production is growing
20 globally. This causes challenges like how to solve increasing energy grid balance problems;
21 how to design and optimize local energy production, consumption and the use of the energy
22 storage; how to cut or shift consumption; and how to operate with more fluctuating energy
23 prices. In the near future this means new businesses and huge global markets to Information
24 and Communication Technology (ICT) solutions for smart grid management in addition
25 to ICT solutions for smart grid adaptable buildings. These challenges are difficult to solve
26 without reliable and scalable forecasting of energy consumption at household, building, block,
27 and district levels. One important piece in the solution of these challenges is to study models
28 to characterize a customer's power consumption with a few parameters.

29 By utilizing Automatic Meter Reading (AMR) data of residential energy consumption,
30 this paper focuses on characterizing and comparing populations of consumers. The aim is
31 to develop efficient and illustrative parameters that allow

- 32 1. comparison and trending of different consumer populations,
- 33 2. communicating all essential elements of consumption, including its volatility, and
- 34 3. easy generation of consumption traces with realistic random variation.

35 The random variation around regular patterns is an essential part of the presented model. In
36 fact, the volatility plays an important role in the control of future low voltage grids, prompt-
37 ing the research beyond the regular consumption patterns. Although the high volatility of
38 the households' energy consumption has been recognized and it is becoming more important
39 in the future grid control, the random variation around consumption patterns has mostly
40 been neglected in modeling.

41 The motivation of this work is to answer the following research questions: i) How to
42 model and parameterize electricity consumption of households with few parameters in such
43 a way that realistic variability in consumption traces can be generated? ii) How to monitor

44 and compare populations of consumers? iii) How to direct towards automatized handling of
45 Big Data by repeated use of autonomous algorithms after initial tuning and validation?

46 The idea of the presented approach is to decompose the consumption into the following
47 components per customer: i) total annual consumption, ii) annual profile, iii) mean week
48 profile and iv) multiplicative random variation, which consists of the difference of the actual
49 consumption time series from the repeating annual and weekly profiles arising from the above
50 components. The authors propose lognormal models for elements i) and iv) and, depending
51 on the context, clustering of the profiles ii) and iii).

52 The main contribution of this work is the identification of multiplicative lognormal noise
53 as a maximally simple way to characterize and monitor the random volatility component by
54 one parameter at minimum. The authors also apply a recently developed grouping (clus-
55 tering) method that is favorable to handle large amounts of data. Consumption clusters
56 and profiles are illustrative and valuable as such, but this approach integrates them into a
57 consumption modeling and monitoring framework as parameters. The overall approach of
58 this paper is holistic, touching many popular problems of energy consumption modeling.

59 There is a vast recent literature on electricity consumption, and the authors bring up
60 only those results from AMR data literature that are closely related to the methodology and
61 ideas of this paper.

62 See McLoughlin et al. (2015) and McLoughlin et al. (2012) for clustering and a review
63 of the same Irish AMR data as utilized in this work. Chicco (2012) presents an overview
64 and performance assessment of the clustering methods for electrical load pattern grouping.
65 A finite mixture model of Gaussian multivariate distributions is introduced in Haben et al.
66 (2016) as an alternative to the popular k-means clustering. This could be an interesting
67 framework to be related to the findings on lognormality, as some of the clustering attributes
68 could benefit from a log-transformation and this work could provide a further attribute
69 describing the random variation. The stability of clustering is studied in Haben et al. (2016)
70 by bootstrapping methods. Clustering of hourly data has been done by utilizing shape

71 dictionaries of consumption patterns and magnitude as a multiplicative factor in Kwac et al.
72 (2014), which is close to the multiplicative decomposition presented here. A multiplicative
73 model in clustering is used in Räsänen et al. (2010) as well, without considering the random
74 variation. Recent developments in clustering include also the clustering of particular time
75 periods and the use of pre-processed load shapes to obtain efficient compression of large data
76 (Kwac et al. 2016; Wang et al. 2017; Haben et al. 2016). Hierarchical methods are used in
77 this context as well, but the focus is directed to grid management, demand response and
78 control, whereas the load prediction for grid control lies outside the scope of this paper.
79 This paper contributes to the methodology of consumption clustering by applying the novel
80 Regular Decomposition clustering method (Reittu et al. 2014; Reittu et al. 2017), which the
81 authors believe to have potential in future needs of clustering, e.g., automated handling of
82 dynamic large data.

83 The proposed hierarchical multiplicative modeling paradigm is motivated by the reported
84 lognormality of energy consumption at various time scales, see (Kuusela et al. 2015; Mutanen
85 et al. 2012; Kwac et al. 2014; Kolter and Ferrera 2011) and the properties of lognormality
86 in Kuusela et al. (2015). For other approaches, see the review Grandjean et al. (2012) of
87 developing consumption traces either by top-down or bottom-up approaches, where the traces
88 in the popular bottom-up approach are obtained by mimicking appliances and generating
89 user and appliance behaviors in various ways.

90 This paper also studies the modeling of the random variation of the annual consumption
91 and the volatility component by relating them to other household characteristics, touching
92 the field of energy consumption survey data analysis. For a review of studies and factors
93 affecting electricity consumption, see Jones et al. (2015), Gouveia and Seixas (2016), and
94 Beckel et al. (2014). Clustering and survey are combined in a recent paper by Gouveia
95 and Seixas (2016). This points also to earlier studies on survey methods in electricity con-
96 sumption. Beckel et al. (2014) extracted 34 features of consumption to reveal household
97 characteristics from the very same data as used in this paper.

98 The outcomes of this paper can be useful to practitioners in various ways. The analysis
99 section provides relatively simple methods for stochastic simulation of the consumption to
100 be utilized, e.g., in populating network models and designing demand response programs
101 as well as in designing new architectural setups, algorithms and decision support tools to
102 utilize distributed energy resources in meeting the demands. Moreover, the authors take a
103 viewpoint of monitoring household energy consumption and propose an intuitive and efficient
104 collection of variables to be measured and monitored. This provides means for electricity
105 distribution companies to trend, compare and predict the consumption. In this framework, it
106 is possible to study both the profiles and their clustering together with the random variation
107 around cluster profiles. The utilized data provide an opportunity to model a household's
108 total consumption without interference of households' own energy production or demand
109 response. Models for the total consumption are needed in, e.g., smart city research. On the
110 other hand, the energy consumption of households is changing in the near future due to the
111 increase in distributed generation and smart devices. This work provides means to observe
112 this change in different scales via trends in monitoring variables.

113 The paper is structured as follows. The research data are summarized in Section 2. In
114 Section 3, the four-layer model is presented and its accuracy is studied. Section 4 is devoted
115 to grouping the annual and weekly profiles by the Regular Decomposition method. The
116 problem of monitoring electricity consumption at a population level is addressed in Section 5
117 by finding suitable consumption monitoring parameters. Possibilities to make inference on
118 electricity consumption based on additional information about households are discussed in
119 Section 6. Finally, the proposed method is validated with unseen data in Section 7. The
120 conclusions are drawn in Section 8.

121 **CHARACTERISTICS OF THE RESIDENTIAL SMART METER DATA AND** 122 **SURVEY**

123 The developed methods are illustrated with the popular dataset of the Irish Smart Meter-
124 ing Trial Archive (2012). The Irish trial took place during 2009 and 2010. The data include

125 smart meter readings at 30 min intervals and participant background data in a survey for-
 126 mat. The data set analyzed in this paper covers 995 households during the 364 first days
 127 (i.e., full weeks) of year 2010. The sample was selected by including all customers heating
 128 their house with electricity (either by central heating or using plug-in heaters) and randomly
 129 picking from the rest homes having uninterrupted records. The authors were interested to
 130 study how the different heating methods are reflected in the electricity consumption traces.
 131 This amounted to the inclusion of 238 homes with electrical heating and 757 homes heated
 132 with other energy sources. Besides the Irish data, the elements of this methodology were
 133 developed with Finnish urban consumer data containing both households and small and
 134 medium-sized enterprises (SMEs).

135 HIERARCHICAL ANALYSIS IN A MULTIPLICATIVE FRAMEWORK

136 The hierarchical analysis of electricity consumption presented in this section will be the
 137 modeling framework through the whole paper. The weeks are indexed by $i = 1, \dots, 52$ and
 138 the half-hour time intervals of a week by $t = 1, \dots, 336$. The electricity consumption C of a
 139 household H in half-hour t of week i is then written as

$$140 \quad C_t^H(i) = W^H \times \frac{y^H(i)}{52} \times \frac{a_t^H}{336} \times \xi_t^H(i), \quad (1)$$

141 where

- 142 W^H = the total annual electricity consumption of the household
- $y^H(i)$ = the weight of week i in the household's annual consumption profile
- a_t^H = the weight of half-hour t in the household's mean week profile
- $\xi_t^H(i)$ = the relative multiplicative variation of consumption around the mean
 profile in week i and time t .

143 The annual and weekly profiles are scaled so that $1/52 \sum_{i=1}^{52} y^H(i) = 1$ and $1/336 \sum_{t=1}^{336} a_t^H =$
 144 1 and, by construction, the irregular variation $\xi_t^H(i)$ has mean 1 as well. Thus, the annual
 145 total consumption is the only element with an absolute magnitude, while the annual and
 146 weekly profiles present the relative distribution of the consumption in time.

147 Energy consumption has been reported to follow lognormal distribution at several time
148 scales: annual scale in Kuusela et al. (2015) and Mutanen et al. (2012), daily scale (with
149 lognormal mixtures) in Kwac et al. (2013), and hourly scale in Chen and Cook (2012) and
150 Mutanen et al. (2012). Kolter and Ferrera (2011) present log-log plots of energy consumption
151 vs. living area, together with the lognormality of the former. In general, many natural
152 phenomena are multiplicative and generate lognormal distributions (Limbert et al. 2001).
153 Multiplication preserves lognormality, which in part suggests the chosen multiplicative model
154 and consumption profile approach.

155 The next step is to analyze these hierarchical elements one by one with the aim to study
156 distributions of variables and suitable models for random variation. In this section, each
157 household H has individual consumption parameters W^H , y^H , a^H , and a model parameter
158 for ξ . For clarity, however, the superscript H will be dropped from the notation in the next
159 section.

160 **Model elements for a single customer**

161 As already mentioned, lognormal modeling of the annual consumption W was studied
162 comprehensively in Kuusela et al. (2015), and the authors adopt it in this paper as well.
163 Besides a mostly good fit with the body of the empirical distribution, lognormal modeling
164 allows heavy tails as well as straightforward transfer of the simple characterization of depen-
165 dencies in multivariate Gaussian models. Figure 1 presents the body and tail fits (inset) of
166 the present data with both lognormal and Weibull distributions. The inset shows that the
167 Weibull distribution underestimates large consumptions considerably. Due to other reported
168 results on lognormality of electricity consumption in various time scales in Kuusela et al.
169 (2015), Mutanen et al. (2012), Chen and Cook (2012), Kwac et al. (2013), and Kolter and
170 Ferrera (2011) as well as the ease of modeling with lognormal distributions, the lognormal
171 model is preferred.

172 Let us then consider the annual and weekly profiles of a household. Recall that in this
173 paper the term 'profile' means a vector with mean one. Define a customer's year profile $y(\cdot)$

174 as the vector of the consumption in each of the 52 weeks divided by the total consumption
 175 W and multiplied by 52 so that the mean of the vector's components is 1. Similarly, define
 176 for each week i the household's profile $\lambda_t(i)$ as the vector of the half-hour consumptions
 177 divided by the total consumption in week i , multiplied by 336. The mean week profile of the
 178 household is then defined as

$$179 \quad a_t = \frac{1}{52} \sum_{i=1}^{52} \lambda_t(i). \quad (2)$$

180 The ratio

$$181 \quad \xi_t(i) = \frac{\lambda_t(i)}{a_t}, \quad i = 1, \dots, 52, \quad t = 1, \dots, 336 \quad (3)$$

182 can now be defined as the multiplicative random variation of the household's energy con-
 183 sumption around its mean profile during week i . Thus, the profiles have been decomposed
 184 multiplicatively as $\lambda_t(i) = a_t \xi_t(i)$. Note also that $\frac{1}{52} \sum_{i=1}^{52} \xi_t(i) = 1$ for every t .

185 Although $\xi_t(i)$ depends on i , it was noticed that most households have rather stable week
 186 profiles in the sense that the process $\xi_t(i)$ retains its character over varying i . Remarkably, the
 187 overall marginal distribution of $\xi_t(\cdot)$ was found to be close to lognormal for most households.
 188 The Kolmogorov-Smirnov distance (maximum deviation between distribution functions) be-
 189 tween each household's random variation distribution and a fitted lognormal distribution is
 190 illustrated in Figure 2. The distance varies in $[0,0.1]$ with mean 0.049, but with few large de-
 191 viations. Although the deviation is big in some individual cases, the marginal distribution of
 192 the multiplicative random variation is mostly well approximated by a lognormal distribution.

193 Since the random variation ξ has mean 1 by construction, its approximating lognormal
 194 distribution is characterized by a single parameter, for example by $\text{Var}(\log(\xi))$. Moreover,
 195 Figure 3 shows that the parameters $\text{Var}(\log(\xi))$ of households are themselves lognormally
 196 distributed.

197 In the following, $\xi_t(i)$ will be modeled by a stationary process with a lognormal marginal
 198 distribution. Before doing this, it is worth of considering the nature of this simplification in
 199 detail. Most households behave qualitatively similarly as the following example (household

29). The left plot of Figure 4 shows the whole process $\log \xi_t(i)$ over a year: a steady random “cloud”. The visual homogeneity is, however, deceptive, because the variance of $(\log \xi_t(\cdot))$ turns out to vary with t with a strong daily pattern. The right plot of Figure 4 presents, for each $t \in \{1, \dots, 336\}$, the empirical variance of $(\log \xi_t(i))_{i=1}^{52}$ (blue curve). The mean week profile of the household is shown for comparison (red curve). Note that the variance is not a monotone function of the mean profile value. Moreover, their shapes differ widely from household to household. However, rough lognormality holds also in this detailed level — the parameter of each lognormal variable then just depends on t , and this picture is very similar for most households. Such a model would, however, be unattractive for practical purposes. We leave now the challenge of more accurate modeling of the $\xi(\cdot)$ processes for the future and look for maximally simple models.

Most mean profiles a_t vary strongly in t (see examples in Section 4), and their marginal distributions can be rather considered as approximately lognormal, with mean one, than approximately Gaussian. An important observation made in this work is that the variation $\xi_t(i)$ depends very weakly on the weekly mean variation a_t . Both time series are close to lognormal, so their logarithms are close to Gaussian, and their dependence is well captured by the respective correlation. The uncorrelatedness of $\log a$ and $\log \xi$ is equivalent to the equality $\text{Var}(\log \lambda) = \text{Var}(\log a) + \text{Var}(\log \xi)$. (The values of $\text{Var}(\log \lambda)$ and $\text{Var}(\log \xi)$ are estimated using all the weeks, and one can use $\log 0 = 0$ when needed.) Figure 5 shows to what extent this holds. The numbers $\text{Var}(\log \lambda)$, $\text{Var}(\log a)$, and $\text{Var}(\log a) + \text{Var}(\log \xi)$ are plotted for all households in the order of increasing $\text{Var}(\log \lambda)$. The first is almost always equal to or a bit smaller than the third, i.e., logarithms of the mean profile a and the multiplicative random variation ξ are slightly negatively correlated. Figure 6 shows these correlations in the same order as the previous figure, and they range between $[-0.2, 0]$, with mean -0.089 .

A study of the temporal behavior of the variation process ξ indicates that, in the mean, the random variation (relative to the household’s mean week profile) that happens at a time

227 point is almost uncorrelated to what happens after 12 hours, but clearly (0.2) positively
 228 correlated to what happens after 24 hours and even after 48 hours again. This is illustrated
 229 in Figure 7.

230 In order to take into account the, albeit small, dependence between ξ and a , as well
 231 as a part of the time correlation, the authors propose modeling the ξ -processes as being
 232 conditioned on the mean week profile process a , and fitting the lag 1 cross-correlations. Note
 233 that although a is non-random, its variance and lag 1 correlation may be computed as for any
 234 time series. By forming for each household H time series ξ_n and a_n over all measurements,
 235 $n = 1, \dots, 336 \times 52$, the empirical covariance matrix $\text{Cov}(\log \xi_n, \log \xi_{n+1}, \log a_n, \log a_{n+1})$
 236 is calculated, where notation n and $n + 1$ refers to studying consecutive measurements. By
 237 averaging all such covariance matrices over all households, the mean covariance matrix

$$\begin{aligned}
 & \text{Mean}(\text{Cov}(\log \xi_n, \log \xi_{n+1}, \log a_n, \log a_{n+1})) & (4) \\
 & = \begin{bmatrix} 0.635 & 0.338 & -0.042 & -0.039 \\ 0.338 & 0.635 & -0.029 & -0.042 \\ -0.042 & -0.029 & 0.391 & 0.357 \\ -0.039 & -0.042 & 0.357 & 0.391 \end{bmatrix}
 \end{aligned}$$

240 is obtained. As $\text{Cov}(\xi_n, a_{n+1})$ and $\text{Cov}(\xi_{n+1}, a_n)$ differ, the process is not invertible in time.

241 **Synthesis of simulated consumption**

242 The mean covariance matrix (4) suggests that the random variation (‘noise’) processes
 243 $\log \xi_t^H(i)$ be almost independent from the mean profiles $\log a_t^H$, whereas the processes $\log \xi_t^H(i)$
 244 have strongly positive lag 1 autocorrelation and differ therefore clearly from white noise.
 245 Note that (4) presents the mean of all household-specific covariance matrices. In order to
 246 assess the significance of the temporal dependence structure of the variation processes, two
 247 sets of simulated traces were generated for each household: one where the true process
 248 $\log \xi_t^H(i)$ was replaced by mean 1 lognormal i.i.d. random variables with the household-

249 specific variance, and one using instead the household-specific empirical covariance matrix
 250 $\text{Cov}(\xi^H, a^H) := \text{Cov}(\log \xi_n^H, \log \xi_{n+1}^H, \log a_n^H, \log a_{n+1}^H)$. Figure 8 compares the variability in
 251 measured and model-generated consumptions at 30 min intervals.

252 The correlated random variation traces produces the same amount of variance for non
 253 electric heating consumers (index range 1 - 757). For heaters, the model overestimates the
 254 variability. The simple i.i.d. model produces larger variability than the one in measured
 255 traces, but the difference is not dramatic, and also this model could be satisfactory for some
 256 purposes. Figure 9 provides details on how well the minimum, median, and maximum of
 257 the weekly consumption maximum are reproduced. The correlated model is slightly better
 258 than the i.i.d. one in predicting the median maximum consumption, while both clearly tend
 259 to produce too large overall maxima. Figure 8 suggests that electric heaters might form a
 260 special group in the modeling.

261 **GROUPING OF ANNUAL AND WEEKLY PROFILES**

262 In the previous section, the consumption traces of households were modeled by the hier-
 263 archical multiplicative model with parameters

$$264 \quad (W^H, y^H, a^H, \text{Cov}(\xi^H, a^H)) \quad \text{or} \quad (W^H, y^H, a^H, \text{Var}(\log \xi^H)), \quad (5)$$

265 where the household-specific profiles y^H and a^H are vectors with lengths 52 and 336, respec-
 266 tively. There is no a priori theoretical model that would generate the observed variety of
 267 annual and weekly mean profiles of a household population. However, the number of model
 268 parameters can be reduced drastically by replacing the individual annual and average week
 269 profiles by mean profiles of relatively homogeneous subgroups of the population. By doing
 270 so, two similarly grouped populations can be compared with each other by comparing the
 271 relative sizes of corresponding groups in the two populations.

272 In order to test the stability of the proposed methodology, the original data were split
 273 into two populations in such a way that both populations had about 24% of heaters. This

274 gives rise to a monitoring development population of 746 households and a test population
 275 of 249 households. The latter will be used in Section 7 to validate the outcome of the present
 276 section.

277 **Grouping of consumption profiles by Regular Decomposition**

278 *The Regular Decomposition method*

279 Clustering algorithms typically divide a data set into groups of elements that are near each
 280 other according to some metric. In contrast, the recently developed Regular Decomposition
 281 method (Reittu et al. 2014; Reittu et al. 2017; Pelillo et al. 2016) aims at a grouping that
 282 is optimal in terms of an information-theoretic criterion, the Minimum Description Length
 283 Principle, Grünwald (2007). Consider a set of customers \mathcal{C} , each having a non-negative time
 284 series $(x_t^{(c)})_{t \in T}$, $c \in \mathcal{C}$. Let \mathcal{P} be a finite partition of \mathcal{C} . As explained in (Reittu et al. 2014;
 285 Reittu et al. 2017), the quantity

$$286 \quad \text{Comp}(x^{(\cdot)}|\mathcal{P}) = \sum_{B \in \mathcal{P}} \sum_{c \in B} \sum_{t \in T} D\left(x_t^{(c)} \left\| \frac{1}{|B|} \sum_{c' \in B} x_t^{(c')}\right.\right), \quad (6)$$

287 where $D(\beta||\alpha) = \alpha - \beta + \beta \log(\beta/\alpha)$ is the Kullback-Leibler divergence between the distribu-
 288 tions $\text{Poisson}(\alpha)$ and $\text{Poisson}(\beta)$, estimates the dominant term of the bit length of a code that
 289 describes the data assuming that the partition \mathcal{P} captures all structure (non-randomness)
 290 present in it (The full code contains also other terms with lesser order of magnitude). For
 291 each positive integer k , the partition

$$292 \quad \mathcal{P}_k^* = \arg \min_{|\mathcal{P}|=k} \text{Comp}(x^{(\cdot)}|\mathcal{P}) \quad (7)$$

293 presents the best grouping into k blocks. Finally, a practically optimal k can be identified as
 294 the smallest k for which the improvement $\text{Comp}(x^{(\cdot)}|\mathcal{P}_k^*) - \text{Comp}(x^{(\cdot)}|\mathcal{P}_{k+1}^*)$ remains below
 295 some small threshold value. Note that popular clustering methods like k-means lack an
 296 inherent principle for the selection of k .

297 It is remarkable that such a grouping can be found in a computationally efficient way.
298 The algorithm presented in Reittu et al. (2014) starts with a random grouping into k blocks
299 and proceeds as a greedy optimization algorithm. As discussed in Reittu et al. (2017),
300 Regular Decomposition has its roots in the mathematics of large structures like graphs and
301 tensors, suggesting a generic applicability of this approach in the separation of structure
302 and randomness in large data. The authors prefer to use the word *group* in the context of
303 Regular Decomposition, as the word *cluster* suggests that the cluster members be close to
304 each other in some metric, which need not always hold.

305 *Grouping of annual profiles*

306 A regular decomposition of the annual profiles suggested six groups denoted by A...
307 F, see Figure 10. The model development data contains about 180 households heating with
308 electricity, but only the group C with 80 members has a large difference between summer and
309 winter consumption. (Recall that the consumption values are scaled, so the profiles show
310 how a household's total consumption spreads throughout the year.) The second largest
311 group B has almost steady consumption throughout the year. The small groups D, E, and
312 F are similar, but D and E show an increase in consumption levels for the third quarter Q3
313 suggesting cooling or other summer time usage. The last three groups are quite small, but
314 the authors wanted to keep these. The analysis in Section 3 showed that the heaters differ
315 to some extent from non-heaters, and the authors hoped to catch the group of heaters by
316 detecting a usage pattern that differs from the majority (the outcome will be examined later
317 in this paper).

318 The obtained profile shapes are remarkably similar to the ones in Gouveia and Seixas
319 (2016), where a grouping was done by Ward's method involving both the pattern and the
320 magnitude of the consumption, contrary to the present method that separates those two. The
321 degree of independence of the total consumption level and profile grouping will be examined
322 in Section 5.

323 *Grouping of average week profiles*

324 Figure 11 illustrates the regular decomposition of the average week profiles and the rich
325 variety in the mean profiles in each group, denoted by a . . . i. Now the optimal number of
326 groups is clearly higher than what was needed for the annual profiles, and there are no very
327 small groups. Most profiles show Mon-Fri vs. Sat-Sun patterns, and these reveal different
328 weekly rhythms of the households' activities.

329 In contrast to McLoughlin et al. (2015) and Kwac et al. (2014), the grouping was done
330 for full weeks instead of individual days. In McLoughlin et al. (2015), the median was used
331 to select a daily profile that a household used most of the time, putting more weight on
332 weekday patterns. Kwac et al. (2014) had a day shape dictionary of 1000 shapes, and they
333 addressed the variability of day shapes by performing an entropy analysis. The groupings in
334 Haben et al. (2016) and Kwac et al. (2016) use particular time periods of day to group the
335 consumption with one European and one US dataset. The key time periods vary somewhat
336 depending on the consumer population.

337 The authors found that households have rather constant weekly rhythms, and the con-
338 sumption evolution through the days of a week is by itself interesting. The authors have
339 also performed an unpublished analysis of urban consumer data containing households and
340 SMEs over 20 districts that illustrated different characteristic weekly rhythms in residential,
341 commercial, and SME industrial districts.

342 **Comparison of measured and groupwise synthesized traces**

343 This section examines at household level the impact of replacing a household's individual
344 annual and weekly profiles by the ones obtained as the average of the profiles within its
345 annual and weekly profile group, respectively. An example is shown in Figure 12 with 30 min
346 consumption traces of a two week period. The measured traces at the top are compared with
347 two alternative synthetic counterparts. The middle row presents the simplest multiplicative
348 model of Section 3 that models the random variation ξ^H by the i.i.d. random variable. The
349 bottom row replaces the individual profiles by the means of their groups. Both synthesized

350 traces are similar to each other as the magnitude of the random variation exceeds the impact
351 of the difference in the profile component. As expected, both models produce larger peak
352 consumption values than the measured ones as illustrated in Figure 9. In addition, the
353 measured traces show more clearly the underlying regular profile shape than the synthesized
354 traces. A shortcoming of the random variation model is seen at the maximum consumption
355 level, and there is also too large variability when the consumption is small or moderate. The
356 authors leave the further tuning of the random variation component model for future work
357 and continue here with the monitoring approach.

358 MONITORING THE PARAMETERS OF A POPULATION

359 The authors propose that a population of energy consumers could be monitored by cal-
360 culating the following variables from the AMR data:

- 361 1. total annual consumptions: the parameters of a lognormal distribution
- 362 2. annual profile: profiles and the frequencies of profile groups
- 363 3. average weekly profile: profiles and the frequencies of profile groups
- 364 4. random variation around the profiles: the parameters of lognormal distributions.

365 This would result in four variables per consumer, i.e., total consumption W , annual profile
366 group, week profile group and a model parameter of random variation, ξ , such as $\text{Var}(\log(\xi))$.
367 In addition to these, there would be $N \times 52$ and $M \times 336$ matrices containing annual and
368 weekly group profile vectors, respectively, with grouping the population into N annual and
369 M weekly groups. By following and comparing these variables, the essential characteristics of
370 residential electricity consumption can be captured. These form a feasible set of monitoring
371 parameters in the following sense: i) they have the power to represent relevant aspects of
372 consumption realistically, ii) the set is minimal and the variables are almost independent from
373 each other (see below), iii) the estimation of parameters is robust, and iv) the comparison
374 of consumer populations is easy. The comparison of populations can be done by comparing
375 lognormal distributions and the frequencies of annual/weekly profiles. It is also easy to

376 generate artificial populations for network models and demand response studies by picking
377 consumer parameters independently from each other.

378 **Independence of the total consumption level and the annual profile group:**
379 Chi square testing of the total consumption (taken with a granularity of 5 MW) and the
380 annual consumption groups shows a dependence between variables due to the three very small
381 groups (that, moreover, have low total consumption levels). When those groups, comprising
382 only 34 members, are removed, the chi square test value becomes 0.84. Thus, for the rest
383 of the data, the annual profile group and the total annual consumption are independent of
384 each other.

385 **Independence of the annual and weekly profile groups:** The annual and weekly
386 profile groups show weak dependence in the model development data (and independence
387 in the test data). By studying the mutual information values between partitions and the
388 expected information between corresponding random partitions, the authors conclude that
389 the annual and weekly profile groups are not informative on each other and can be considered
390 as independent from each other.

391 **RELATING HOUSEHOLD CHARACTERISTICS TO CONSUMPTION**

392 **PARAMETERS**

393 This section takes advantage of the associated survey data in order to model the total
394 annual consumption and the random variation. Relating the household characteristics to the
395 consumption is not necessary for monitoring purposes, but such models would allow deeper
396 understanding and offer more possibilities in the generation of new realistic consumption
397 populations. The authors attempted to find a profile classifier based on the household
398 characteristics. However, no valuable linkage was found, even for central heaters. This
399 outcome is in line with Gouveia and Seixas (2016), McLoughlin et al. (2012), and McLoughlin
400 et al. (2015). A low correlation between energy usage behavior and geodemographics is also
401 reported in Haben et al. (2013).

402 **Stochastic models for the total annual consumption and the random variation**
 403 **parameter**

404 The number of persons, the number of rooms, and the home floor area increase a house-
 405 hold’s energy consumption (Gouveia and Seixas 2016; Jones et al. 2015). These will be
 406 related to the total annual consumption and the amount of random variation. The authors
 407 have applied these characteristics successfully in earlier research with Finnish and Irish data
 408 to model the total annual consumption (Kuusela et al. 2015). Moreover, using such data
 409 is practical as it is typically available, and it is close to housing district planning data as
 410 well. Data mining methods applied to the Irish data in Beckel et al. (2014) were successful
 411 in inferring the occupancy, the number of persons and the number of appliances and, with
 412 some difficulties, the floor area and the number of bedrooms from 34 energy consumption
 413 features derived from the dataset.

414 This paper utilizes the multivariate lognormal model and the notation from Kuusela
 415 et al. (2015) to derive multivariate lognormal distributions for the vectors (P, F, B, W) and
 416 (P, F, B, V) , where P =number of persons + 0.5, F =home floor area, and B =number of
 417 bedrooms + 0.5, and W = total consumption in MWh, $V = \text{Var}(\log(\xi))$ (the addition of
 418 0.5 to P and B is only for plotting purposes). The multivariate lognormal distribution is
 419 parameterized by $\boldsymbol{\mu}$, the vector of mean values of the log-transformed variables, and Γ , the
 420 covariance matrix of the log-transformed variables. The estimated model parameter $\boldsymbol{\mu}$ equals
 421 $(1.185, 5.014, 1.4254, 2.172)$ and the parameter Γ equals

$$\begin{bmatrix} 0.190, & 0.053, & 0.032, & 0.110 \\ 0.053, & 0.160, & 0.053, & 0.083 \\ 0.032, & 0.053, & 0.047, & 0.046 \\ 0.110, & 0.083, & 0.046, & 0.280 \end{bmatrix}, \tag{8}$$

423 for the vector (P, F, B, W) . The respective parameters for the (P, F, B, V) -vector are

424 (1.185, 5.014, 1.425, -0.435) and

$$425 \begin{bmatrix} 0.190, & 0.053, & 0.032, & -0.015 \\ 0.053, & 0.160, & 0.053, & -0.045 \\ 0.032, & 0.053, & 0.047, & -0.021 \\ -0.015, & -0.045, & -0.021, & 0.211 \end{bmatrix}. \quad (9)$$

426 The estimation results are listed for the two estimations as the interest will be to estimate
427 W or V , given P, F , and B .

428 Figure 13 presents the marginal densities of multivariate lognormal fits to the target
429 variables at the top row. The lognormal distribution fits very well to W and V . Lognormal
430 fits are rather good for P and F as well, but the variable B is skewed to the opposite direction
431 in comparison to the other variables. The granularity and the concept of a bedroom might
432 be a bit problematic, see the discussion in Kuusela et al. (2015).

433 VALIDATION WITH THE TEST POPULATION

434 This section studies i) the stability of the grouping of the annual and weekly consump-
435 tion profiles and ii) the ability to predict the total annual consumption and the random
436 variation parameter by household characteristics. Also, the predicted consumption traces
437 are compared with the measured ones.

438 Stability of annual and weekly profiles

439 In the Regular Decomposition method the number of groups as well as the annual and
440 weekly profile vectors were fixed to those obtained from the model development data. Then
441 the same classification algorithm was run to group the annual and average week profiles from
442 the validation data.

443 This grouping with a fixed scheme works well also for the new data; the previously fixed
444 profiles and the averages of profiles among group members are very close to each other. In
445 the four largest annual groups, the fixed schema provides a very good match. Naturally, one

446 should not include groups of insufficient size in population monitoring.

447 Figure 14 illustrates the largest difference in weekly profiles. The differences in profiles
448 are associated with the group size and hence with the averaging over member profiles. Since
449 even the profile pair with the largest difference captures well the essential consumption
450 pattern, the authors conclude that grouping the unseen validation data with fixed weekly
451 profile function works well and allows to compare customer populations by recording the
452 frequencies of profiles in the population.

453 In this validation data, the annual and weekly groups are independent of each other (chi
454 square independence test value 0.13).

455 *Grouping with fixed vs. free profiles*

456 What results if only the number of annual and weekly clusters is fixed, and the cluster
457 profiles are let to be optimal for the test data? This kind of analysis provides information
458 on the goodness of grouping with fixed profiles. Firstly, one needs to verify that the group
459 profiles resulting from optimization are close to the fixed profiles. It is also interesting how
460 the consumers form the groups. This question is examined with the weekly grouping, where
461 all the groups have substantial sizes.

462 It turns out that 64% of the test data is grouped so that there is a very close profile
463 from the fixed development data group profile set. Overall, the new group average profiles
464 are quite similar to the fixed profiles (although less smooth due to the smaller number of
465 samples in the averaging). However, it is not easy to identify a mapping to the whole data
466 set that takes a grouping with fixed profiles to the grouping with free profiles. An interesting
467 observation is that the consumer groups do not remain unchanged when the profiles are let
468 to be free. The variation in households' individual average week profiles is still large and
469 hence the memberships of the groups are not always obvious. However, the resulting group
470 average profiles are rather stable. Thus, one should not follow the group membership labels
471 of individual consumers in time, but what kind of groups the consumers form. 78% of the
472 validation data is covered by the five largest groups, and it is rather easy to find a mapping

473 between the fixed group mean profiles (from the model development data) and the new group
474 mean profiles (from the validation data). The closest profile can be chosen unequivocally in
475 four cases, and the remaining one has a few rather close profile candidates. The best profile
476 matches are illustrated in Figure 15.

477 **Prediction of the annual consumption and the random variation**

478 In this section, the total annual consumption W and the random variation parameter V
479 are estimated by conditioning each on the household size P , the home floor area F , and the
480 number of bedrooms B . The estimators are the conditional expectation of W given (P, F, B)
481 and that of V given (P, F, B) , derived in Kuusela et al. (2015). The conditional distribution
482 of W given (P, F, B) , denoted as $W|(P, F, B)$, is lognormally distributed, and similarly for
483 V . Formulas for the expectations and the variances of $W|(P, F, B)$ and $V|(P, F, B)$ can be
484 written by equations (2) and (3) of Kuusela et al. (2015). It turns out that the conditional
485 expected value cannot predict the target variables accurately at the household level. This
486 is due to the large variability of households: the estimator is the expected value of the
487 conditional consumption. Instead, the models can reproduce a similar random variation
488 in the target values as that existing in the test population (see the discussion in Kuusela
489 et al. (2015)). However, the selection of consumers for this paper results to a worse model
490 than the one analyzed more deeply in Kuusela et al. (2015). For each validation observation
491 (P, F, B, W) , the conditional distribution $W|(P, F, B)$ and its 95% confidence interval was
492 formed. In this validation sample, 18% of the W values were outside of the 95% confidence
493 intervals compared to less than 5% in a population of the same Irish data utilized in Kuusela
494 et al. (2015). Note that the conditional distribution is a function of (P, F, B) so that the
495 confidence intervals are also functions of these variables.

496 The random variation component is studied with the group profiles obtained by fixing the
497 annual and weekly profiles to those obtained from the model development data. In less than
498 4% of the observed test data, the random parameter values are outside the 95% confidence
499 interval of the random value parameter estimator. The pair of curves in Figure 16 illustrates

500 the distributions of the observed random variation parameter values and the corresponding
501 model-generated values obtained by picking 50 samples from the conditional distribution
502 given the household characteristics of each validation data consumer, i.e., $V|(P, F, R)$. The
503 observed random variation parameter values tend to be larger than the ones generated by
504 the developed model, although the difference is not huge. However, it will be visible in the
505 model-generated consumption traces shown in Figure 17, where the model tends to predict
506 a smaller random variation than the observed variation. One possible reason could be that
507 by the grouping with predefined annual and weekly profiles, the random component includes
508 an impact of the non-optimal group profiles in addition to the pure random variation. Thus,
509 the most realistic random variation scheme should use the households' individual profiles.

510 CONCLUSIONS

511 This work contributed to the field of electricity consumption modeling and monitoring
512 by analyzing a multiplicative modeling framework consisting of i) total annual consumption,
513 ii) annual consumption profile, iii) average weekly consumption profile, and iv) random
514 variation around the repeated mean consumption profiles. The variation of consumption is
515 a natural element in the model and very easy to monitor in this framework. This modeling
516 intuition stemmed from the lognormality of the electricity consumption. Section 3 showed
517 that the model was able to sufficiently capture the amount of random variation around the
518 repeated consumption patterns, and the generated consumption traces accurately reproduce
519 the minimum and the median of a consumer's weekly consumption maxima. However, the
520 random variation model would benefit from further tuning at low and, in particular, at peak
521 consumption levels.

522 Then the interest was turned towards monitoring a population of electricity consumers
523 and the properties of the proposed monitoring parameters. For that purpose, the recently
524 developed Regular Decomposition method was utilized to group the annual and weekly
525 profiles. It turned out that the monitoring parameters were essentially independent from
526 each other. The validation showed good stability of the groups. The authors propose to

527 direct research interest towards the random variation around regular patterns as the amount
528 of randomness exceeds small differences in profiles. The grouping of profiles would benefit
529 from efficient methods to handle dynamic large data.

530 The data provide an opportunity to model the households' total electricity consumption
531 as household energy systems were rare in Ireland during the trial period. When the house-
532 holds' energy production and smart energy systems will become common, it will be very
533 difficult to assess the actual energy consumption of a household, as the energy companies
534 only see the amount of energy required to meet the total consumption. The rapid evolution
535 in household energy equipment and the offered energy products as well as the tariffs also have
536 an impact on the data collected by the energy companies, and the modeling of households'
537 total consumption will become increasingly difficult.

538 The developed household consumption model offers a relatively simple method to simulate
539 the stochastic variation of electricity consumption to populate network models or to design
540 new architectural setups, algorithms, and decision support tools to utilize distributed energy
541 resources in meeting the demands.

542 **References**

- 543 Archive, I. S. S. D. (2012). "CER Smart Metering Project. [http://www.ucd.ie/issda/
544 data/commissionforenergyregulationcer/](http://www.ucd.ie/issda/data/commissionforenergyregulationcer/).
- 545 Beckel, C., Sadamori, L., Staake, T., and Santini, S. (2014). "Revealing household charac-
546 teristics from smart meter data." *Energy*, 78, 397 – 410.
- 547 Chen, C. and Cook, D. J. (2012). "Behavior-based home energy prediction." *Intelligent En-
548 vironments (IE)*, 2012 8th International Conference on, IEEE, 57–63.
- 549 Chicco, G. (2012). "Overview and performance assessment of the clustering methods for
550 electrical load pattern grouping." *Energy*, 42(1), 68–80.
- 551 Gouveia, J. P. and Seixas, J. (2016). "Unraveling electricity consumption profiles in house-
552 holds through clusters: Combining smart meters and door-to-door surveys." *Energy and
553 Buildings*, 116, 666–676.

554 Grandjean, A., Adnot, J., and Binet, G. (2012). “A review and an analysis of the residential
555 electric load curve models.” *Renewable and Sustainable Energy Reviews*, 16(9), 6539–6565.

556 Grünwald, P. (2007). *Minimum Description Length Principle*. The MIT Press.

557 Haben, S., Rove, M., Greethm, D., Grindord, P., Holderbaum, W., Potter, B., and Single-
558 ton, C. (2013). “Mathematical solutions for electricity networks in a low carbon future.”
559 *22nd International Conference on Electricity Distribution (CIRED), Stockholm, 10-13*
560 *June 2013*.

561 Haben, S., Singleton, C., and Grindrod, P. (2016). “Analysis and clustering of residential
562 customers energy behavioral demand using smart meter data.” *IEEE Transactions on*
563 *Smart Grid*, 7(1), 136–144.

564 Jones, R. V., Fuertes, A., and Lomas, K. J. (2015). “The socio-economic, dwelling and ap-
565 pliance related factors affecting electricity consumption in domestic buildings.” *Renewable*
566 *and Sustainable Energy Reviews*, 43, 901 – 917.

567 Kolter, J. Z. and Ferrera, J. J. (2011). “A large-scale study on predicting and contextualizing
568 building energy usage.” *Proceedings of the Twenty-Fifth AAAI Conference on Artificial*
569 *Intelligence, 7-11 August 2011, San Francisco, California, USA. AAAI Press, 2011*.

570 Kuusela, P., Norros, I., Weiss, R., and Sorasalmi, T. (2015). “Practical lognormal framework
571 for household energy consumption modeling.” *Energy and Buildings*, 108, 223–235.

572 Kwac, J., Flora, J., and Rajagopal, R. (2014). “Household energy consumption segmentation
573 using hourly data.” *Smart Grid, IEEE Transactions on*, 5(1), 420–430.

574 Kwac, J., Flora, J., and Rajagopal, R. (2016). “Lifestyle segmentation based on energy
575 consumption data.” *IEEE Transactions on Smart Grids*, PP(99).

576 Kwac, J., Tan, C.-W., Sintov, N., Flora, J., and Rajagopal, R. (2013). “Utility customer
577 segmentation based on smart meter data: Empirical study.” *Smart Grid Communications*
578 *(SmartGridComm), 2013 IEEE International Conference on*, 720–725 (Oct).

579 Limbert, E., Stahel, W., and Abbt, M. (2001). “Log-normal distributions across the sciences:
580 keys and clues.” *BioScience*, 51(5), 341 – 352.

581 McLoughlin, F., Duffy, A., and Conlon, M. (2012). “Characterising domestic electricity
582 consumption patterns by dwelling and occupant socio-economic variables: An Irish case
583 study.” *Energy and Buildings*, 48, 240–248.

584 McLoughlin, F., Duffy, A., and Conlon, M. (2015). “A clustering approach to domestic
585 electricity load profile characterisation using smart metering data.” *Applied Energy*, 141,
586 190 – 199.

587 Mutanen, T., Sorasalmi, T., and Weiss, R. (2012). “Electricity usage type selection and
588 model validation based on electricity usage measurements.” *Industrial Conference on Data
589 Mining 2012. IbaI Publishing (2012)* (July).

590 Pelillo, M., Elezi, I., and Fiorucci, M. (2016). “Revealing structure of large graphs: Sze-
591 merédi’s regularity lemma and its use in pattern recognition.” *Pattern Recognition Letters*.

592 Räsänen, T., Voukantsis, D., Niska, H., Karatzas, K., and Kolehmainen, M. (2010). “Data-
593 based method for creating electricity use load profiles using large amount of customer-
594 specific hourly measured electricity use data.” *Applied Energy*, 87(11), 3538–3545.

595 Reittu, H., Bazsó, F., and Norros, I. (2017). “Regular Decomposition: an information and
596 graph theoretic approach to stochastic block models arXiv:1704.07114.

597 Reittu, H., Weiss, R., and Bazso, F. (2014). “Regular decomposition of multivariate time
598 series and other matrices.” *Structural, Syntactic, and Statistical Pattern Recognition*, P.
599 Fränti, G. Brown, M. Loog, F. Escolano, and M. Pelillo, eds., number 8621 in LNCS,
600 Berlin-Heidelberg, Springer-Verlag, 424–433.

601 Wang, Y., Chen, Q., Kang, C., Xia, Q., and Luo, M. (2017). “Sparse and redundant
602 representation-based smart meter data compression and pattern extraction.” *IEEE Trans-
603 actions in Power Systems*, 32(3).

604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630

List of Figures

1	Histogram of the total consumption distribution (blue) and its fits with a lognormal distribution (red) and a Weibull distribution (green), inset shows the tail probability fit in \log_{10} scale.	27
2	Kolmogorov-Smirnov distances (maximum deviation between distribution functions) between each household's random variation distribution and a fitted lognormal distribution.	28
3	Histogram of the log-variances of the random variation processes of all customers, and its fit with a lognormal distribution (red).	29
4	(a): The process $\log \xi_t(i)$ of household 29 over a whole year. (b): The empirical variance of $(\log \xi_t(i))_{i=1}^{52}$ for $t \in \{1, \dots, 336\}$ (blue curve), and the mean week profile of household 29 (red curve).	30
5	The numbers $\text{Var}(\log \lambda)$ (black, thick), $\text{Var}(\log a)$ (blue), and $\text{Var}(\log a) + \text{Var}(\log \xi)$ (red) for each household, plotted in the order of increasing $\text{Var}(\log \lambda)$	31
6	The correlations between $\log a$ and $\log \xi$ for each household, plotted in the order of increasing $\text{Var}(\log \lambda)$	32
7	The average, over all households, of the logarithmic autocorrelations of the random variation processes for half-hour lags $1, \dots, 96$	33
8	Variances of the normalized consumption processes of all customers, presented as a cumulative plot. Blue: the true values. Green: synthetic traces with i.i.d. random variation. Red: synthetic traces with correlated random variation. Customers with numbers 758-995 use electrical heating.	34
9	The smallest, the largest, and the median values of weekly maxima of the normalized consumption processes of all customers, presented as a cumulative plot. Colors and indexing as in Fig. 8.	35
10	Annual consumption profiles A . . . F of the model development data with group size in the parenthesis.	36

631	11	Group mean profiles and sizes in the grouping of weekly average profiles in the model development data.	37
632			
633	12	Comparison of observed,(a), and two modeled, (b) and (d), traces as well as the non-random components of the two models, (c) and (e). The non-random components consist of annual and weekly profiles. The figures (b) and (c) utilize individual profiles of the selected customer whereas figures (d) and (e) illustrate results by utilizing the group mean annual and weekly profiles. . .	38
634			
635			
636			
637			
638	13	Marginals of multivariate lognormal modeling. Figure (a) illustrates the annual consumption, W , in MWh and (b) the variance around consumptions patterns, V . These are the target variables. The remaining figures illustrate the household characteristics persons in household, P in figure (c), home area in m^2 , F in (d), and the number of bedrooms, B in (e). These household characteristics are related to each of the target variables in turn.	39
639			
640			
641			
642			
643			
644	14	Illustration of the largest difference between the fixed week profile (blue) and the mean of group profile (orange) in the validation data set. This group had 16 members.	40
645			
646			
647	15	The best profile matches between weekly profiles of the model development data (blue) and new group mean profiles (red) from free classification of households' average weeks into nine groups in the validation data.	41
648			
649			
650	16	The distributions of the observed (blue) and model generated (red) random variation parameter values in the validation data.	42
651			
652	17	An example of comparison of consumption traces in the validation data. Profile figure illustrates the non-random component from consumer's annual and weekly classification profiles.	43
653			
654			

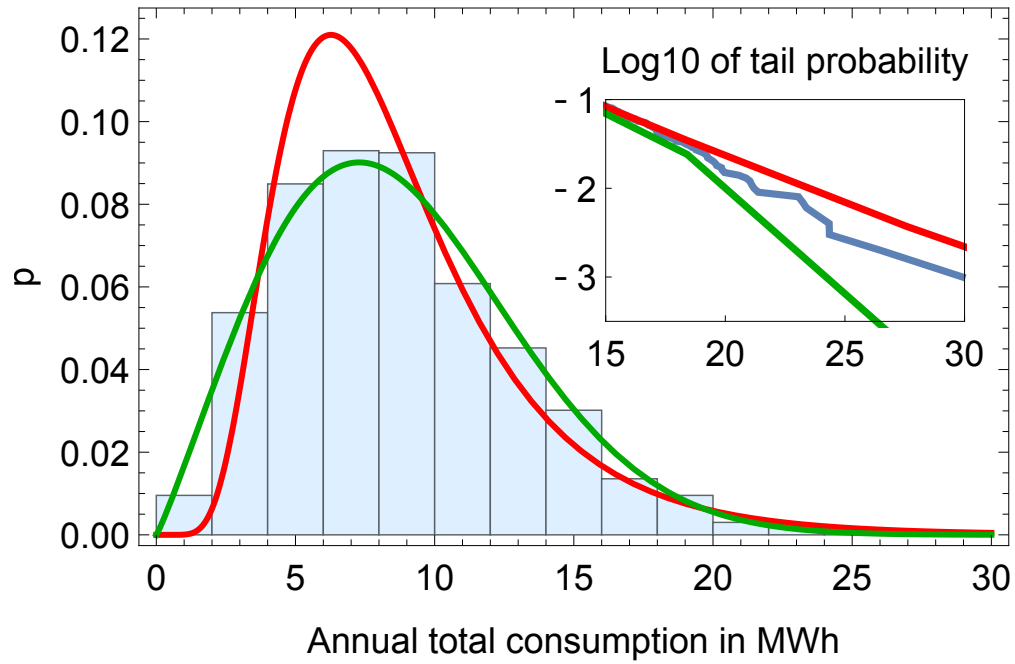


Fig. 1. Histogram of the total consumption distribution (blue) and its fits with a lognormal distribution (red) and a Weibull distribution (green), inset shows the tail probability fit in \log_{10} scale.

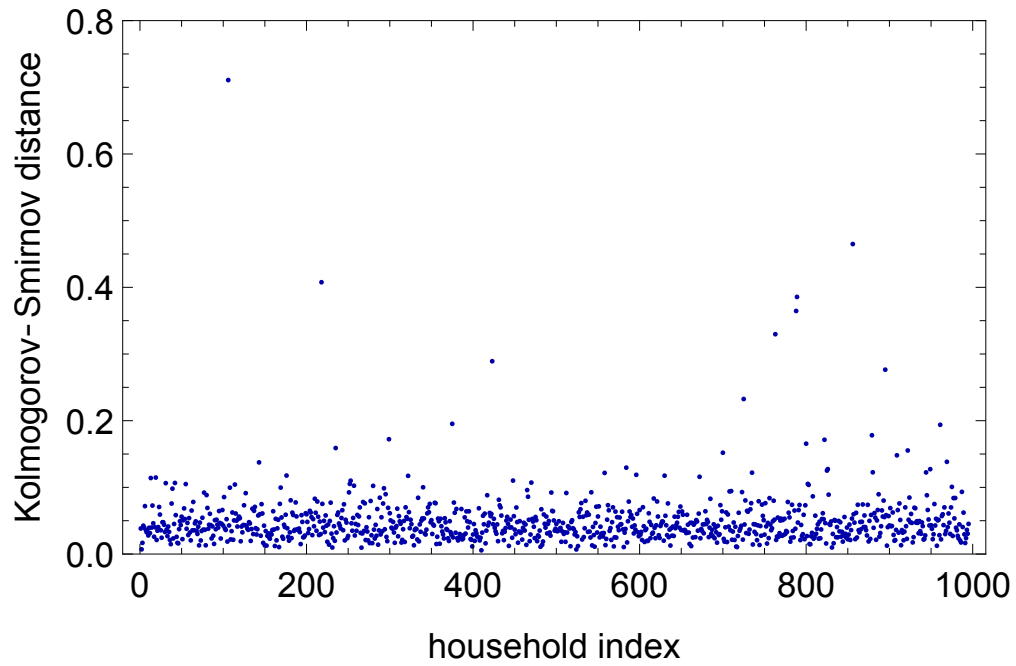


Fig. 2. Kolmogorov-Smirnov distances (maximum deviation between distribution functions) between each household's random variation distribution and a fitted log-normal distribution.

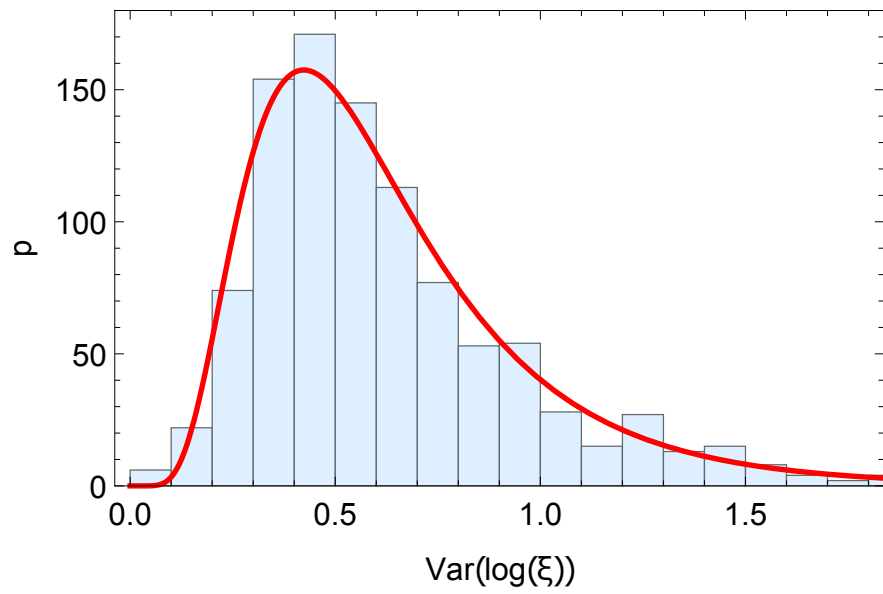


Fig. 3. Histogram of the log-variances of the random variation processes of all customers, and its fit with a lognormal distribution (red).

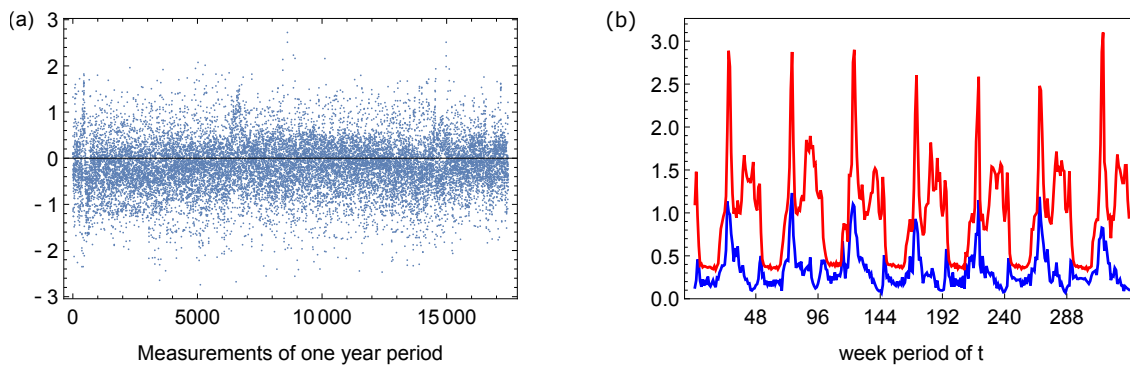


Fig. 4. (a): The process $\log \xi_t(i)$ of household 29 over a whole year. (b): The empirical variance of $(\log \xi_t(i))_{i=1}^{52}$ for $t \in \{1, \dots, 336\}$ (blue curve), and the mean week profile of household 29 (red curve).

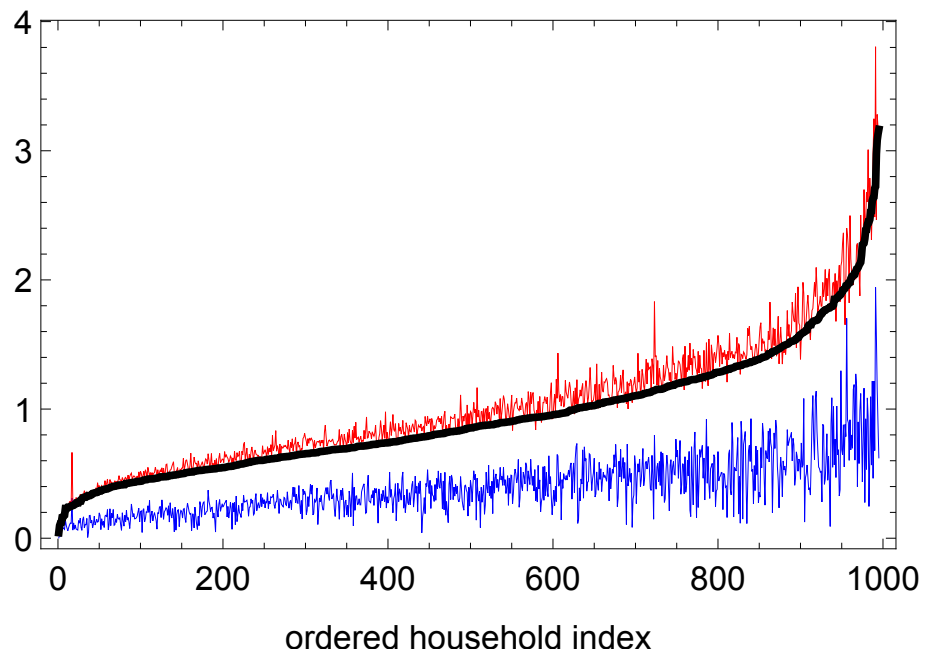


Fig. 5. The numbers $\text{Var}(\log \lambda)$ (black, thick), $\text{Var}(\log a)$ (blue), and $\text{Var}(\log a) + \text{Var}(\log \xi)$ (red) for each household, plotted in the order of increasing $\text{Var}(\log \lambda)$.

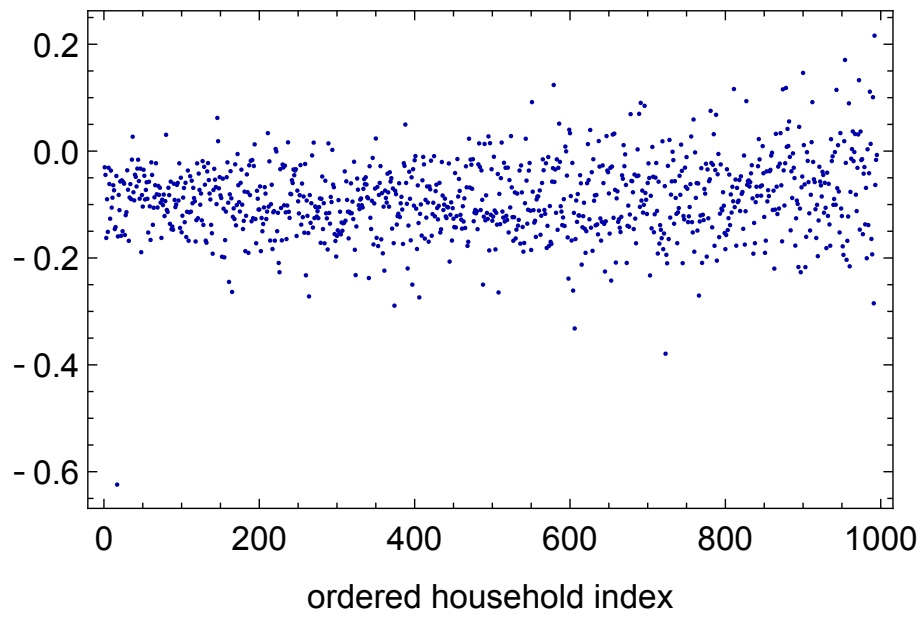


Fig. 6. The correlations between $\log a$ and $\log \xi$ for each household, plotted in the order of increasing $\text{Var}(\log \lambda)$.

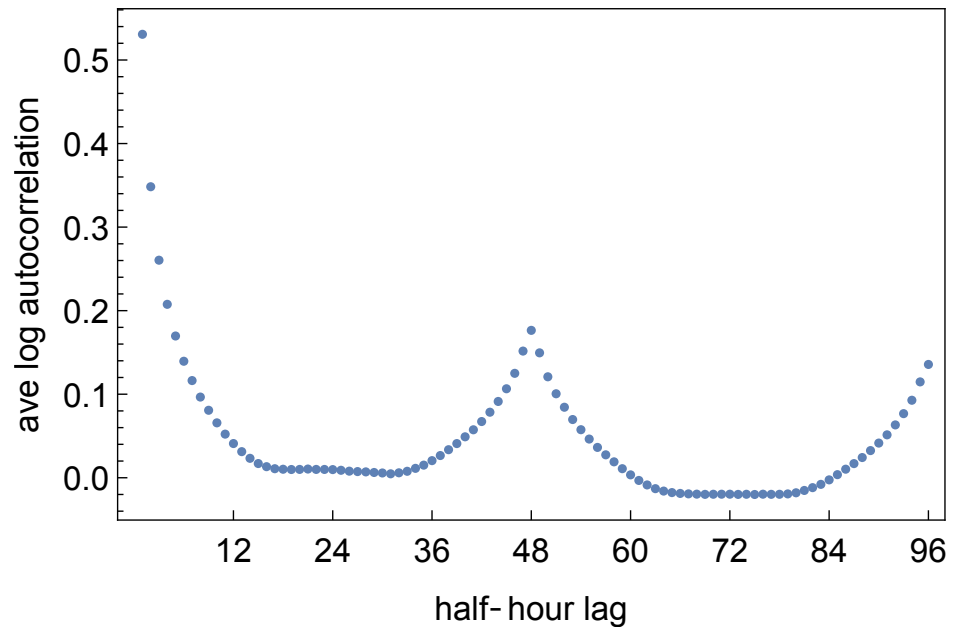


Fig. 7. The average, over all households, of the logarithmic autocorrelations of the random variation processes for half-hour lags $1, \dots, 96$.

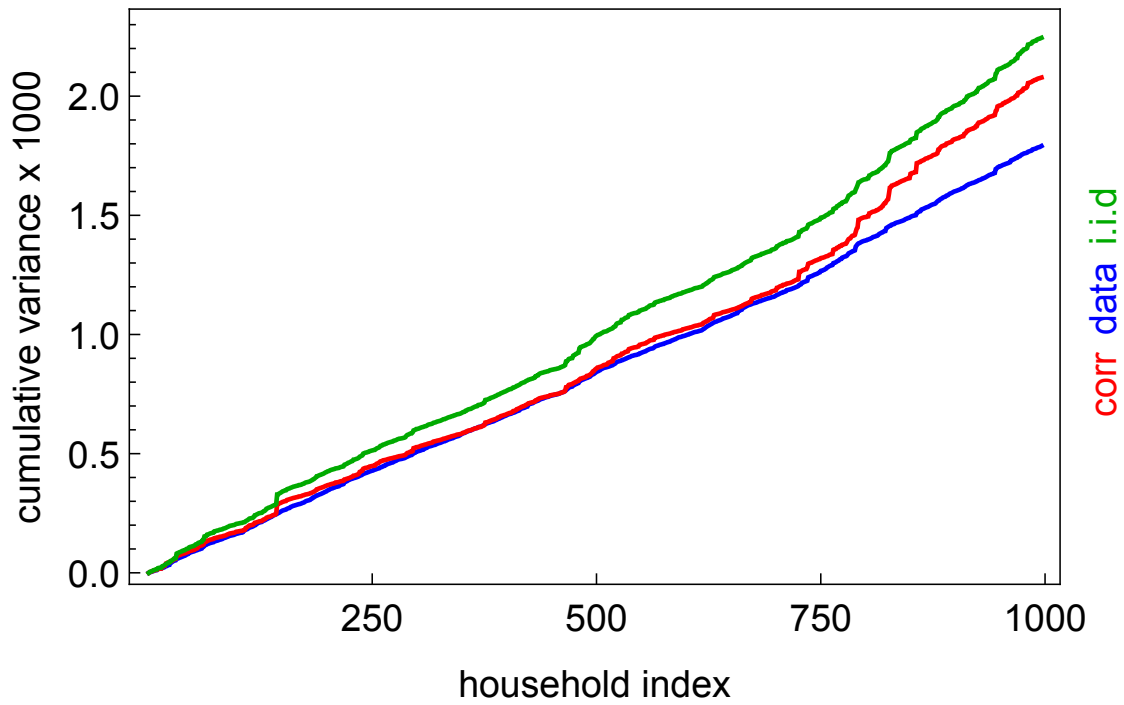


Fig. 8. Variances of the normalized consumption processes of all customers, presented as a cumulative plot. Blue: the true values. Green: synthetic traces with i.i.d. random variation. Red: synthetic traces with correlated random variation. Customers with numbers 758-995 use electrical heating.

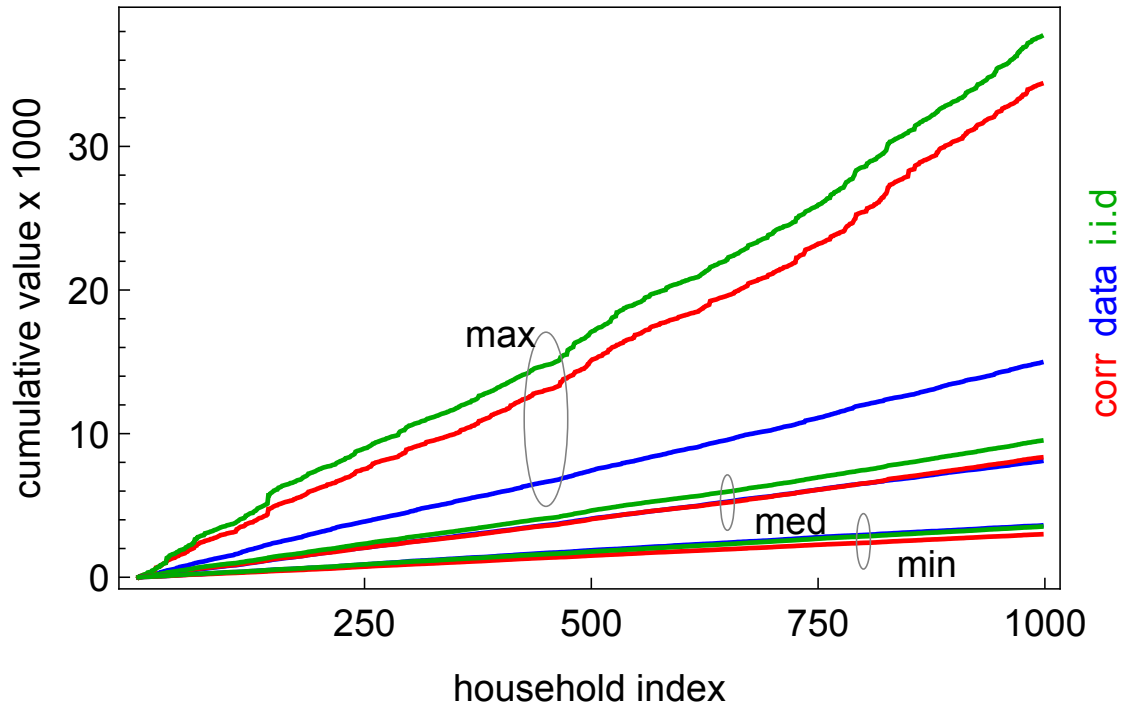


Fig. 9. The smallest, the largest, and the median values of weekly maxima of the normalized consumption processes of all customers, presented as a cumulative plot. Colors and indexing as in Fig. 8.

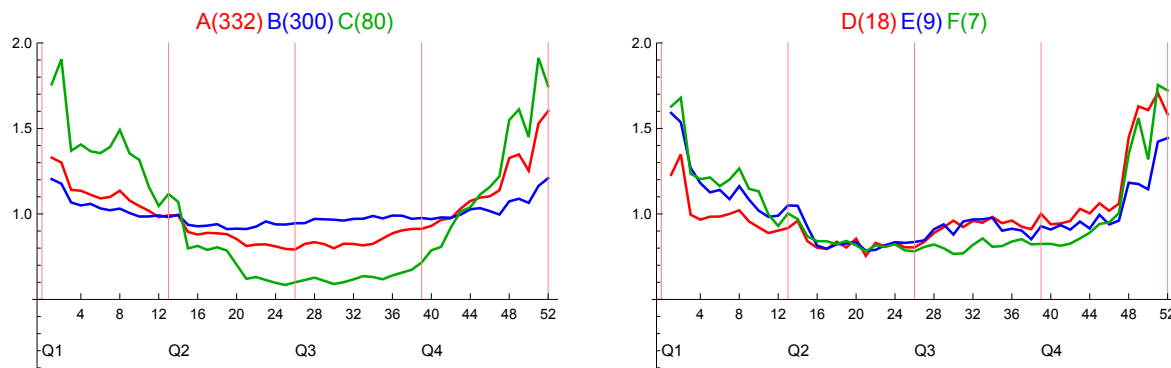


Fig. 10. Annual consumption profiles A... F of the model development data with group size in the parenthesis.

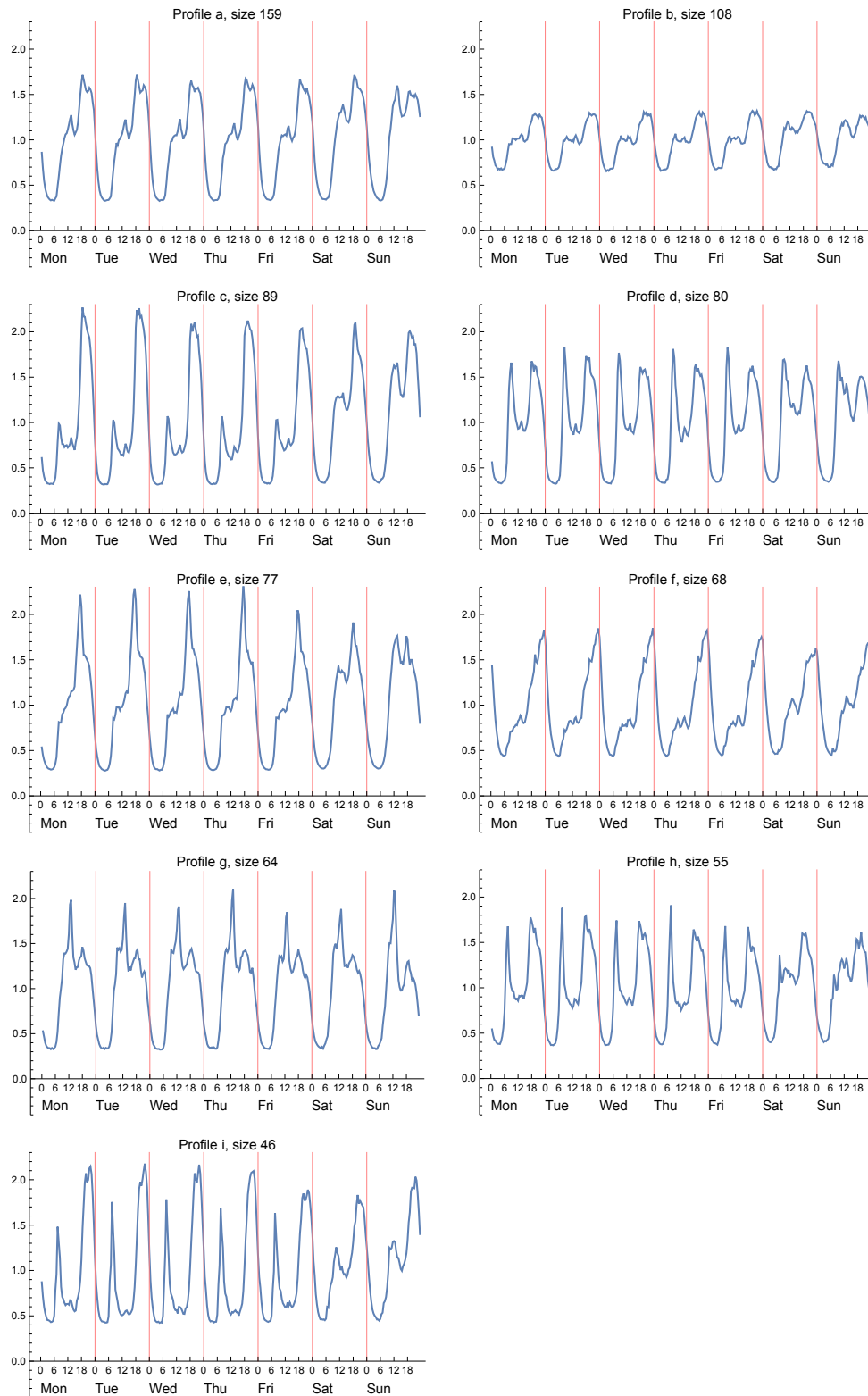


Fig. 11. Group mean profiles and sizes in the grouping of weekly average profiles in the model development data.

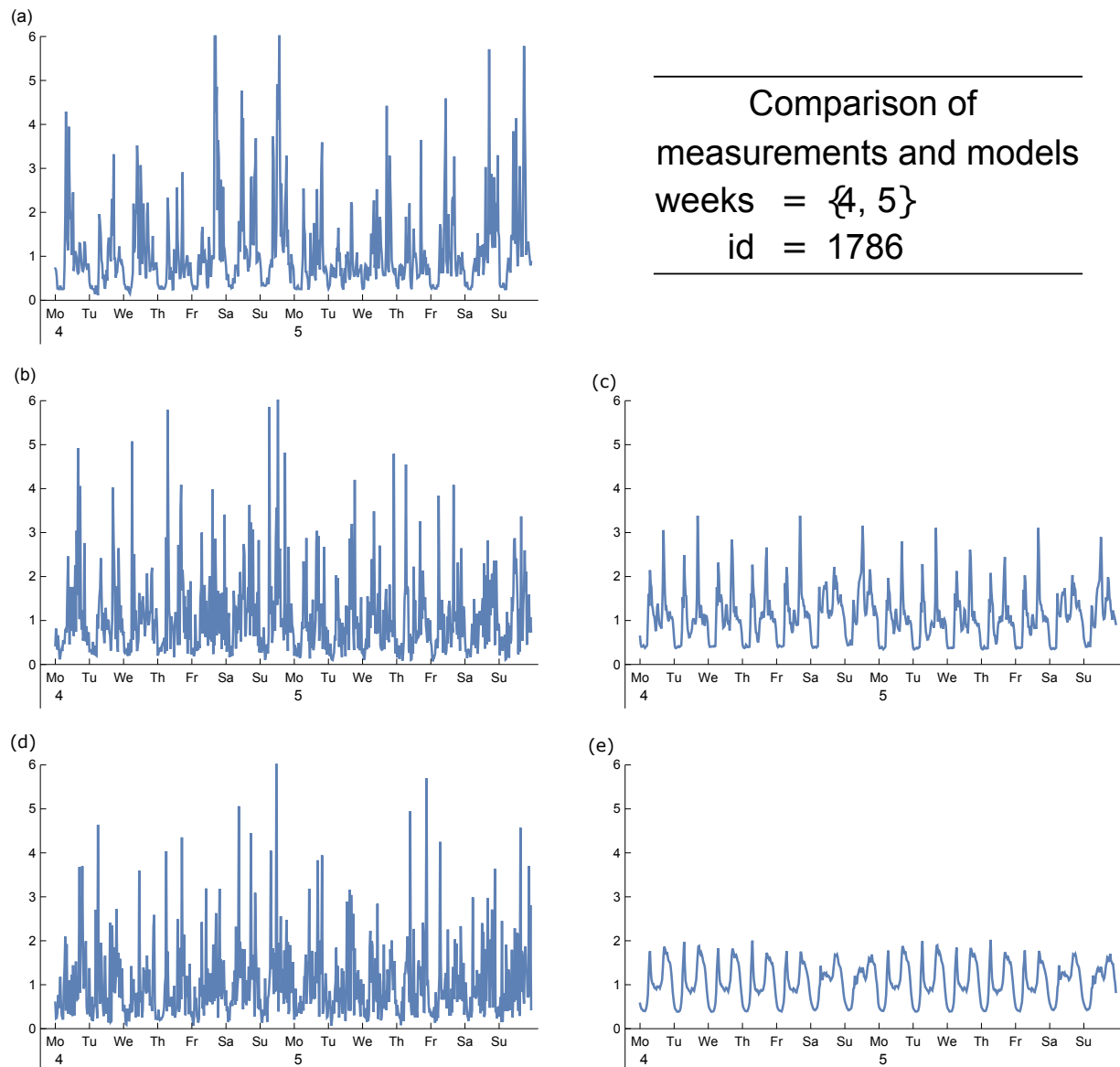


Fig. 12. Comparison of observed, (a), and two modeled, (b) and (d), traces as well as the non-random components of the two models, (c) and (e). The non-random components consist of annual and weekly profiles. The figures (b) and (c) utilize individual profiles of the selected customer whereas figures (d) and (e) illustrate results by utilizing the group mean annual and weekly profiles.

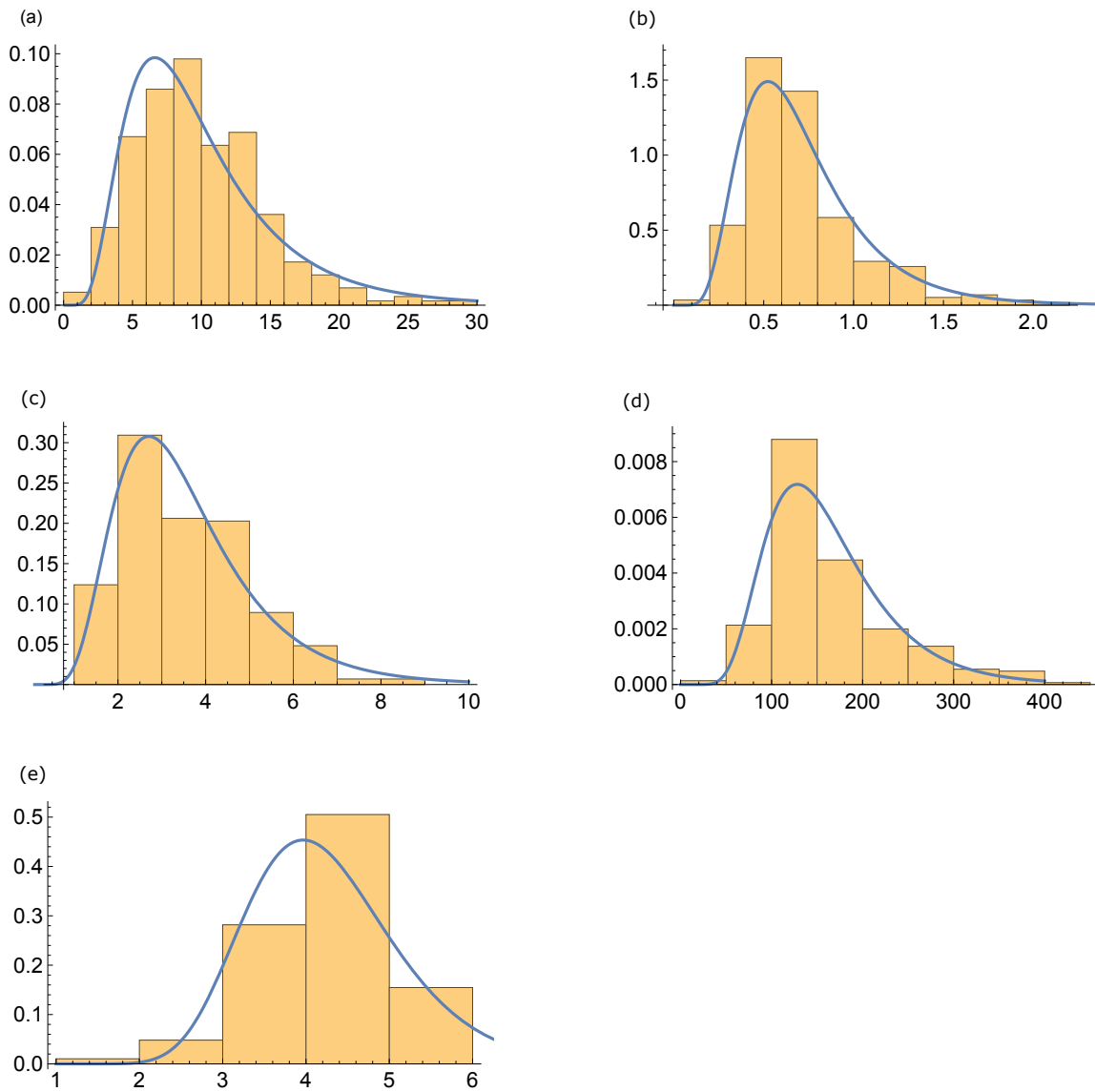


Fig. 13. Marginals of multivariate lognormal modeling. Figure (a) illustrates the annual consumption, W , in MWh and (b) the variance around consumption patterns, V . These are the target variables. The remaining figures illustrate the household characteristics persons in household, P in figure (c), home area in m^2 , F in (d), and the number of bedrooms, B in (e). These household characteristics are related to each of the target variables in turn.

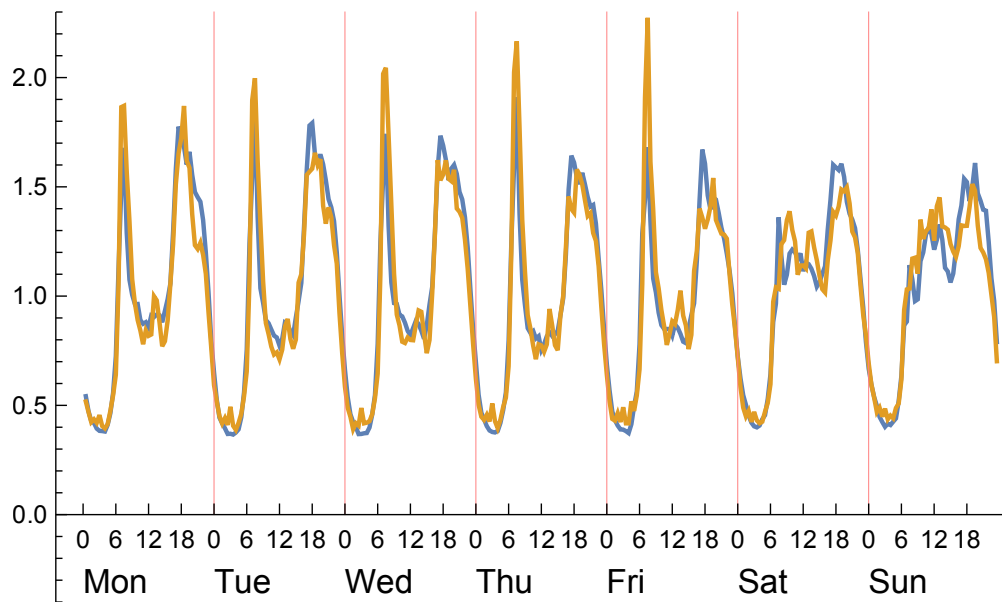


Fig. 14. Illustration of the largest difference between the fixed week profile (blue) and the mean of group profile (orange) in the validation data set. This group had 16 members.

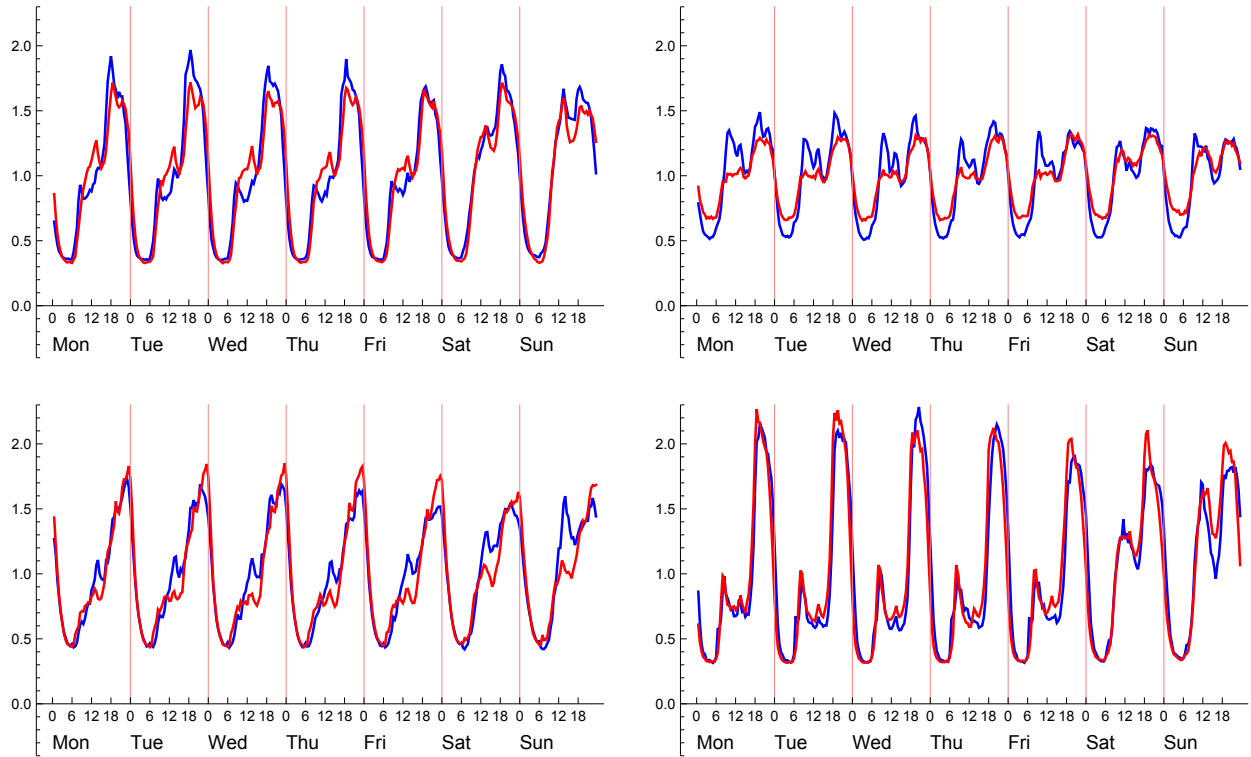


Fig. 15. The best profile matches between weekly profiles of the model development data (blue) and new group mean profiles (red) from free classification of households' average weeks into nine groups in the validation data.

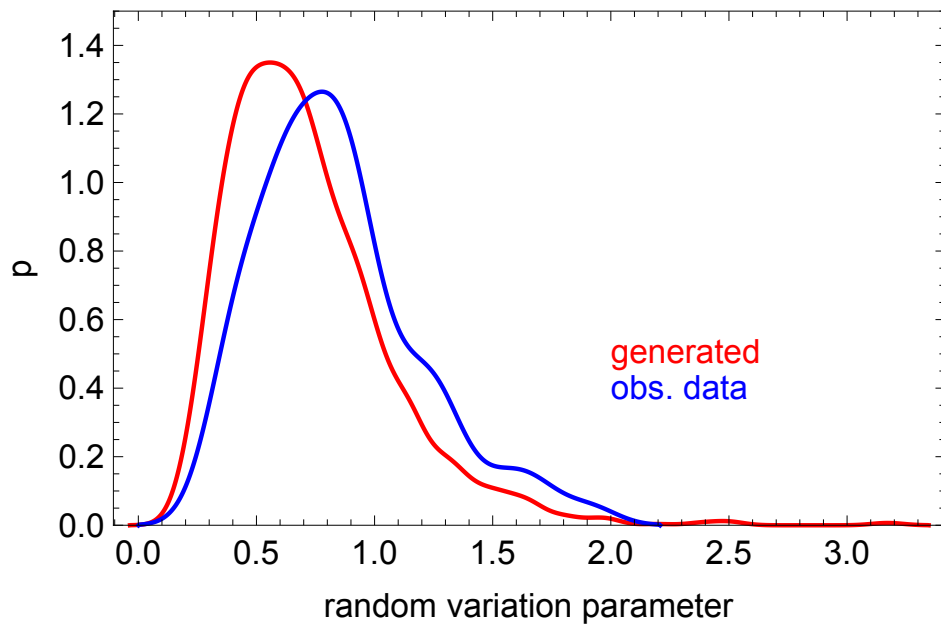


Fig. 16. The distributions of the observed (blue) and model generated (red) random variation parameter values in the validation data.

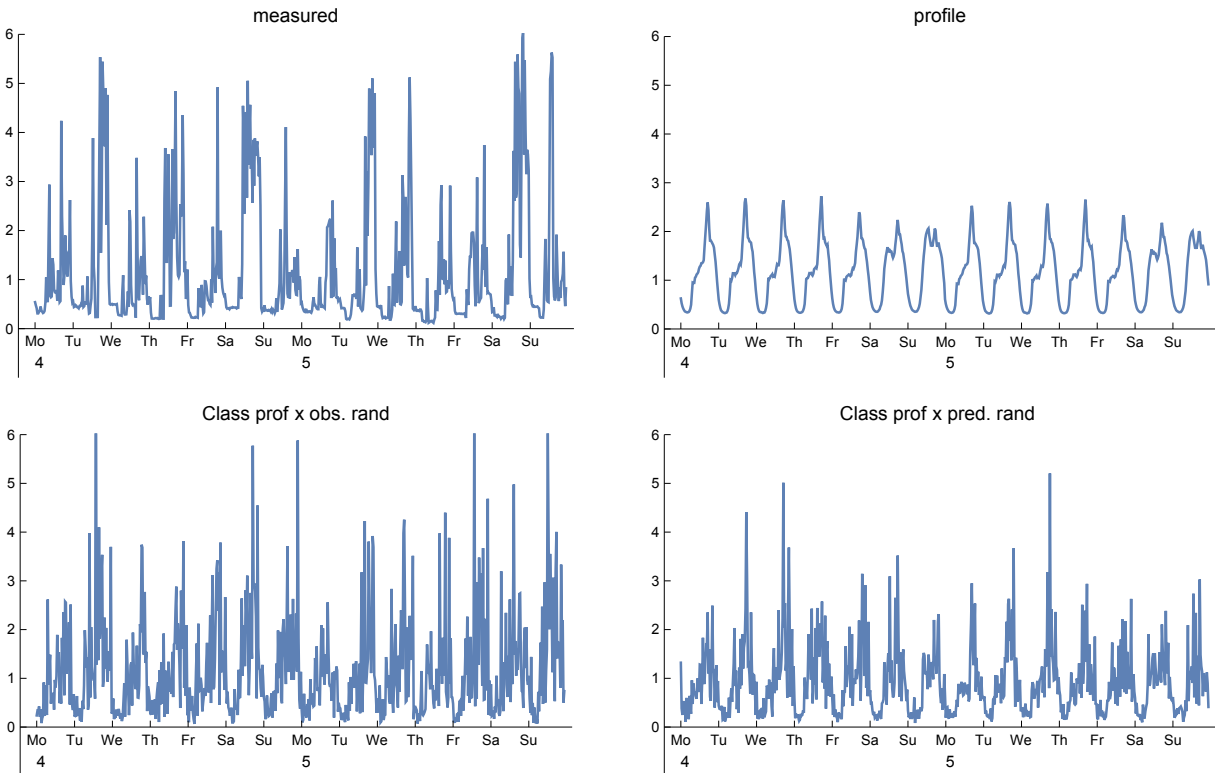


Fig. 17. An example of comparison of consumption traces in the validation data. Profile figure illustrates the non-random component from consumer's annual and weekly classification profiles.