

Households Electricity Consumption Analysis with Data Mining Techniques

Usman Ali, Concettina Buccella, Carlo Cecati

Department of Information Engineering, Computer Science and Mathematics

University of L'Aquila, L'Aquila, Italy

email: mail.usmanali@gmail.com, concettina.buccella@univaq.it, carlo.cecatti@univaq.it

Abstract—Smart Grid improves the electricity grid infrastructure by introducing new powerful communication system between consumer and supplier. Implementation of smart meters increases the availability of detail level of consumer electricity load profile data. To improve and efficient planning and development of this new power system, a primary challenge is to analyze the electricity consumption data. To analyze the energy consumption or achieve our objective we choose the best analytic process is data mining techniques including exploratory data analysis and preprocessing, frequent patterns mining and associations, classification /characterization, clustering and outlier deduction. In this paper, we use these techniques and apply on two different public available datasets. Explain and evaluate which techniques is use full for the better understanding of electricity load profile consumption data.

I. INTRODUCTION

A smart grid is a modernized electrical grid and to better planning and development of this system we must need the electricity consumption analysis is a primary challenge. With the provision of sustainable energy and deregulation in energy market contributed to interest in this field. In smart grid, electricity smart meters are used to store the energy consumption and related information at a different interval and granularity i.e. per second, per minutes and hourly basis. All the energy consumption information collected from smart meters from the different specific geographical area and transmitted in real-time to a central location through a modern intelligent communication system. This centralized power grid can be analyzed for further process. In new power grid system availability of accurate, reliable and updated information on energy consumption helpful for demand forecasting, prediction, energy conservation demand side management, demand response and similar activities to active real-time data-driven analytic operations [4].

To maximize the efficiency of supply process of electricity, and more improvement of the new grid the market provides more targeted and complicated tariff offers for electricity customers. Introduction of smart meters will allow greatly increases in the analysis of a consumer's electricity usage pattern that's also provide the ability to make customized offers on electricity pricing. These analytic also help in demand response or demand side management to minimize electricity usage during peak hour timing or to increase efficiencies in the electricity utility supply chain [5].

Based on consumers of electricity, the usage of electricity is different for each and every consumer resulting different

load shapes due to different consumer behavior, weather situations, different weekdays and time also matters. The primary objective of our load profiling task over a set of users having different recorded attributes and their half hourly demand over a period of month is to build models which approximate their load shapes for certain subsets of customers and self-reliance estimates for those load shapes, for different consumer weather conditions, times of year, and days of the week.

Our secondary objective includes producing models and visualization which helps us to promote understanding underlying structures which may be in load profile data, and to identify the relative importance and interactions between various useful variables. The third objective is understanding of electrical usage patterns within a household is necessary which house level device usage may impact of overall energy usage reduction during different peak times of the day. If we use the smart plugs that give accurate devices level consumption to provide a pattern of total load demand. That would help for the supplier side generation and transmission [5]. To analysis the energy consumption and to achieve our objective we choose the best computational technique is data mining. In data mining, we analyze different data techniques to the problem of extracting meaningful knowledge from the noisy and large database [15]. Our goal is the ability to adapt to the local characteristics of the data is supposed to be the most important feature of data mining which is applied to the electricity load consumption related to various factors such as weather, time and customer characteristics.

In past few years, energy analytic is emerging and lots of research on electricity consumption analysis like consumer segmentation, characterization, predictions and knowledge extraction from smart meter had been done [16]. The data mining techniques mostly used classification, clustering of electricity demand patterns and cluster analysis of smart metering data [8], [6], [13]. However, mostly research on the smart meter consumption analysis, but in this paper we target the consumer (household) and utilities (microgrid) level analysis. In consumer level, our main contribution is analysis not only base on the smart meter but cover all household appliances used in each home. In utility level analysis, we cover the whole microgrid at per minutes granularity.

Keeping objective in above, we use two datasets first we use REDD: A Public Data Set for Energy Disaggregation Research [11]. The second dataset is Smart* microgrid that having load profile more than 400 houses [2]. Our goals are to

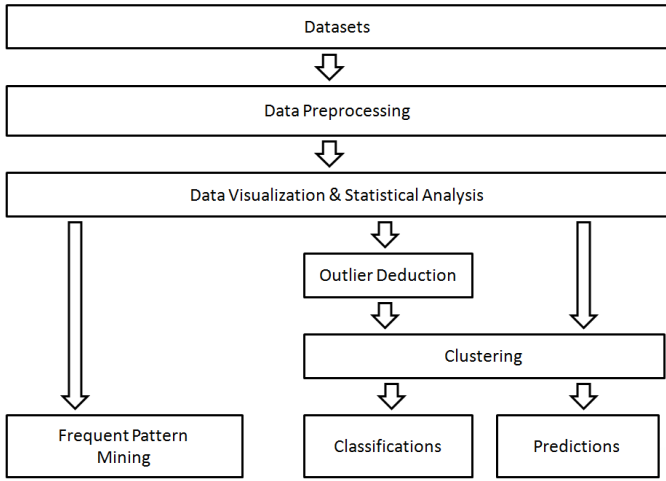


Fig. 1. Process of Data Mining Techniques Apply on datasets

apply following data mining techniques on these datasets like exploratory data analysis and preprocessing ,data visualization techniques, frequent patterns and associations,classification /characterization ,clustering and outlier deduction [14] [9].

II. DATA MINING TECHNIQUES

In this paper, we use different data mining techniques for analysis of electricity load profile. Figure 1 shows the process how we apply different data mining techniques on data and detail discussion on these in the following section.

A. Data Preprocessing

Data in the real world is raw format or incomplete that lacking new attribute values and certain interesting attributes , or sometimes data having only aggregate data like in electricity load profile and sometimes data having noisy containing errors or outliers. Also inconsistent containing discrepancies in the data in the form of codes or names. Major preprocess techniques are follows:

1) *Data cleaning*: The data is filled with any missing values, smooth noisiness in the data, identify or try to remove outliers, and resolve inconsistencies in the data.

2) *Data integration*: In this process the integration of multiple databases, data cubes, or files into a single or useable format.

3) *Data transformation*: In this phase data further process for normalization and aggregation of data.

4) *Data reduction*: In this process reduced representation in volume but produces the same or similar analytical results.

5) *Data discretization*: This process also part of reduction but with particular importance and in the form of ranges ,or binning especially for numerical data.

B. Statistical Analysis

Statistical analysis techniques use in this paper for basic summary and understanding of dataset. This can be done through five number summary like maximum, minimum ,1st

Quartile, 3rd Quartile, Median and Mean. This statistical analysis tells the basic idea how the dataset is useful or not for further experimentation.

C. Frequent Patterns and Associations

In data mining, the most common task of finding the frequent pattern in large-scale databases or datasets is very important and for past few years these techniques have been studied in large scale. First proposed by Agrawal, Imielinski, and Swami in the context of frequent itemsets and association rule mining [1]. The frequent pattern algorithm we use is FP-Growth. The FP-Growth Algorithm, proposed by Han [7], is an efficient and scalable algorithm for mining the complete set of frequent patterns by pattern fragment growth that us using an extended prefix-tree structure for saving the compressed and important information about frequent patterns by using the frequent pattern tree (FP-tree).

D. Clustering

Clustering is a collection of same group similar or data objects or in other groups. Its also finding the dissimilar to the data objects in other groups. In cluster analysis, the main objective is to find similarities between data objects with the help of specific characteristics found in the data and grouping these similar data objects into clusters.

1) *Clustering Algorithms*: In this paper, we use partitioning base clustering method. Partitioning a database D of n objects into a set of k clusters, such that the sum of squared distances is minimized (where c_i is the centroid or medoid of cluster C_i). We use two algorithms K-mean and K-medoids.

K-means (MacQueen's, Lloyd's): Each cluster is represented by the center of the cluster. K-means is an unsupervised partitional classification algorithm, which requires the exact information of the number of clusters in order to operate [12].

K-medoids clustering selects the most centrally located data points within clusters as cluster centers called medoids. The algorithm employed is based on Partitioning Around Medoids (PAM) method (Kaufman and Rousseeuw, 1987) with modifications [10]. Initial medoids are computed by K-means++ centroid initialization algorithm. Cluster medoids are updated by choosing the most centric data point in a then-currently existing cluster, instead of computing a swapping cost for each data point and crosschecking as proposed by the original algorithm.

2) *Clustering Distance Measure*: Clustering algorithms make heavy use of distance between two vectors to compare data points and clusters. Several distance measures exist, but we use Euclidean. Euclidean distance between two Q-dimensional vectors is computed as Equation (1) .

$$d_e(V_i, Z_j) = \|V_i, Z_j\|^2 = \sqrt{\sum_{q=1}^Q (V_i, Z_j)^2} \quad (1)$$

3) *Clustering Prediction*: We also use the clustering algorithm for the prediction of estimated load profile. This can be done using Map Clustering on Labels.

4) *Clustering Evaluation*: The quality of a clustering algorithm can be demonstrated via evaluating generated clusters internally by validity indices that provide means of comparison among distinct clustering results and present the infrastructure for automatically selecting the optimum number of clusters to be generated on a particular data set. Most common measure Sum of Squared Error (SSE). The error is the distance to the nearest cluster for the each points. To calculate the SSE, we square these errors and sum them. You can see in given Equation (2). Where x is a data point in cluster C_i and m_i is the representative point for cluster C_i . Given two clusters, we can choose the one with the smallest error.

$$SSE = \sum_{i=1}^K \sum_{x \in C_i} dist^2(m_i, x) \quad (2)$$

We use Sum of Square Error for evaluation of what will use for the optimum value of k cluster. For the evaluation of prediction results use accuracy, recall, precision and F-measure.

E. Data Classification

Classification is also known as pattern recognition, discrimination, or supervised learning. There is a lot of Classification approaches, including the use of decision trees and rule induction; density estimation; and artificial neural networks. A decision tree is a tree-like graph structure or model. It also looks like an inverted tree because it grows downwards and it has root at the top. The representation of the data provides the user more advantage compared with other methods that are being meaningful and easy to understand or interpret. The goal is to create a classification model that use for prediction.

F. Outlier Deduction

Outlier detection is a crucial and interesting research problem in data mining that goal to find objects that look like a dissimilar, exceptional information and inconsistent with respect to the majority of the data. This technique also part of data preprocessing. In this paper, we use three major outlier deduction techniques.

1) *Distance Based*: In Distance-based outlier detection, an object p is an outlier if its neighborhood does not have enough other points. For each object p , examine the number of other objects in the r -neighborhood of p , where r is a user-specified distance threshold. An object p is an outlier if most (taking — as a fraction threshold) of the objects in D are far away from p , i.e., not in the r -neighborhood of q Equation (3).

$$\frac{||\{q | dist(p, q) < r\}||}{||D||} \leq \alpha \quad (3)$$

An object p is a $DB(r, \alpha)$ outlier if Equivalently, one can check the distance between p and its k -th nearest neighbor q , where $k = \lceil \alpha ||D|| \rceil$. p is an outlier if $dist(p, q) > r$.

2) *Density Based*: In Density-based outlier detection, an object o is an outlier if its density is relatively much lower than that of its neighbors. Use the relative density of an object against its neighbors as the indicator of the degree of the object being outliers.

3) *Local Outlier Factor (LOF)*: Local outlier factor (LOF) is an outlier algorithm that was proposed by Markus M. Breunig in 2000 for finding anomalous data points by calculating the local deviation of an each given data point with respect to its nearest or neighbors [3]. LOF of an object p is the average of the ratio of local reachability of p and those of p 's k -nearest neighbors. The lower the local reachability density of p , and the higher the local reachability density of the k NN of p , the higher LOF. This can be computed as given Equation (4):

$$LOF_k(p) = \frac{\sum_{o \in N_k(p)} \frac{lrd_k(o)}{lrd_k(p)}}{||N_k(p)||} \quad (4)$$

III. EXPERIMENTS AND RESULTS

We use two different datasets for the experimentation of data mining techniques. First, REDD: A Public Data Set for Energy Disaggregation Research [11]. A freely available data set containing detailed power usage information from several homes, which is aimed at furthering research on energy disaggregation (the task of determining the component appliance contributions from an aggregated electricity signal). The data set, containing several weeks of power data of device level for 6 different homes, and high-frequency current/voltage data for the main power supply of two of these homes. The second datasets is the Smart*: An Open Data Set and Tools for Enabling Research in Sustainable Homes. This dataset having Electrical load profile data over a single 24-hour period from 443 unique homes [2].

TABLE I. DEVICES AND MONITORS USE IN REDD DATASET

House	Monitors	Device Categories
1	20	Electronics, Lighting, Refrigerator, Disposal, Dishwasher, Furnace, Washer Dryer, Smoke Alarms, Bathroom GFI, Kitchen Outlets, Microwave
2	11	Lighting, Refrigerator, Disposal, Dishwasher, Washer Dryer, Kitchen Outlets, Microwave, Stove
3	22	Lighting, Refrigerator, Dishwasher, Washer Dryer, Bathroom GFI, Kitchen Outlets, Oven, Microwave, Electric Heat, Stove
4	20	Lighting, Dishwasher, Furnace, Washer Dryer, Smoke Alarms, Bathroom GFI, Kitchen Outlets, Stove, Disposal, Air Conditioning
5	26	Electronics, Lighting, Refrigerator, Disposal, Dishwasher, Furnace, Washer Dryer, Bathroom GFI, Kitchen Outlets, Microwave, Electric Heat, Outdoor Outlets
6	17	Kitchen Outlets, Washer Dryer, Stove, Electronics, Bathroom GFI, Refrigerator, Dishwasher, Electric Heat, Lighting, Air Conditioning

A. Data Preprocessing

The REDD Dataset required a lot of preprocessing. Figure 2 shows how preprocessing perform on datasets. But the microgrid datasets required less processing.

Following steps we did during the preprocessing of datasets:

1) *Data cleaning*: Change the Time Format UTC integer to readable date format. Correct inconsistent data by sorting the time series data in proper format.

2) *Data integration*: We Integrate the dis-aggregated data into one CSV format file. All the devices level data merge into one file of each house.

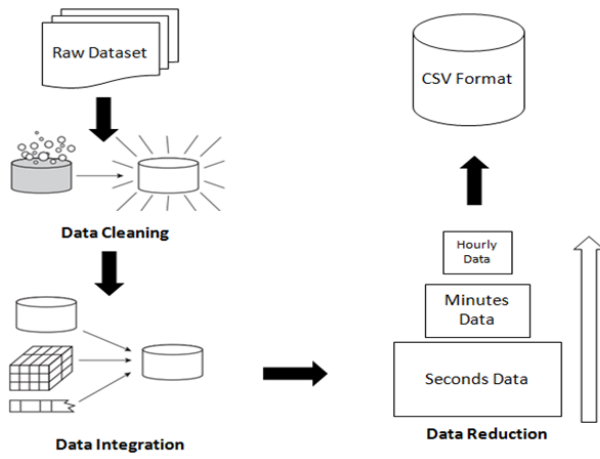


Fig. 2. Data Preprocessing

3) *Data transformation*: Attribute construction: replacing or adding new attributes inferred by existing attributes. Introduce new Attributes like :

Weekdays/Weekend

Morning (6 AM to 10 AM), Noon (10 AM to 5 PM) , Evening (5 PM to 8 PM) , Night (8 PM to 6 AM)

Peak Hour (6 AM to 10 AM or 5 PM to 8 PM) ,Off-Peak Hour

Days of Week (Monday,Tuesday...).

4) *Data reduction*: In our case, we did data reduction on the basis of the time base. The data set is available in seconds and we have to convert into hourly data for understanding or analysis. so we first convert the dataset to minutes and then hours.

TABLE II. CONSUMPTION ANALYSIS OF REDD DATASET

	House 1	House 2	House 3	House 4	House 5	House 6	Total
Weeks							
Weekdays	139.7	35.6	147.9	127.4	39.9	105.7	596.2
Weekend	97.1	15.3	30.9	31.3	1.7	34.4	210.6
Hours							
Off Peak	169.8	34.5	137.8	113.9	31.5	94.3	581.8
Peak	66.9	16.4	41.0	44.8	10.0	45.8	225.0
Week Days							
Sunday	54.1	9.1	17.5	18.0	1.7	15.4	115.8
Monday	39.1	8.2	34.0	24.3	10.0	20.3	135.9
Tuesday	39.9	7.4	27.9	23.2	26.5	22.3	147.3
Wednesday	20.6	7.3	42.7	27.5	2.3	26.1	126.4
Thursday	17.6	7.2	21.6	31.1	1.2	23.9	102.6
Friday	22.5	5.5	21.7	21.3	0.0	13.2	84.2
Saturday	43.0	6.2	13.4	13.3	0.0	19.0	94.8
Periods							
Evening	25.6	7.3	20.8	16.8	2.5	10.3	83.3
Morning	23.5	5.3	12.1	17.1	6.3	25.5	89.8
Night	108.8	21.9	107.4	80.2	22.6	75.7	416.5
Noon	78.8	16.3	38.6	44.7	10.2	28.5	217.2
Total	236.7	50.9	178.8	158.7	41.6	140.1	806.8

B. Data Understanding

Data understanding is done through some statistical or visualization of graphs. Table II shows the detail view of data:

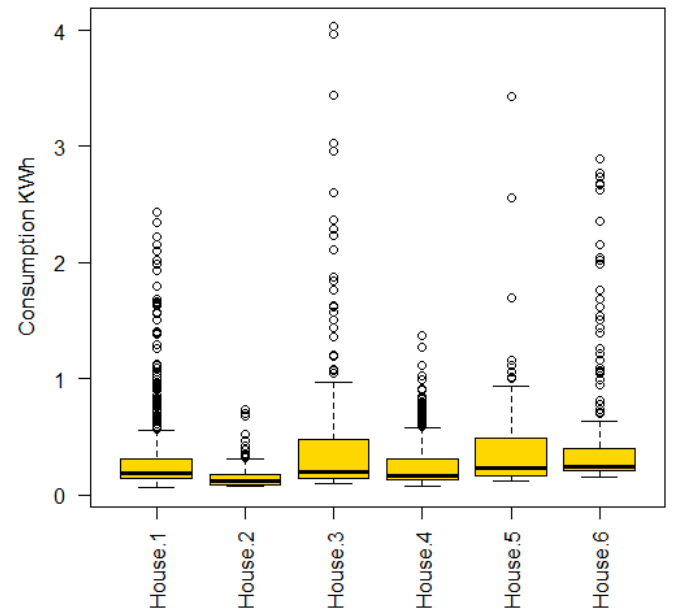


Fig. 3. Boxplot of REDD Datasets Houses

1) *Weeks Analysis*: The weekly analysis shows that the House 1 having more consumption on a weekend than other houses and overall 26% consumption on a weekend.

2) *Days Analysis*: Days of week analysis shows that the House 1 having more consumption on Sunday or Saturday than other houses. House 4 having more consumption on Tuesday and overall all the days having the similar percentage of consumption.

3) *Periods Analysis*: Periods in days analysis shows that mostly consumption on night period overall 52% of the total is cover in night timing. Second most consumption on the noon period is 27% of all the houses.

4) *Hours Analysis*: The hourly analysis shows that on these houses 28% of consumption between the peak hour.

5) *Total Consumption Analysis*: Figure 3 boxplot shows the five number summary of all the 6 houses and we can see that House 3 having maximum load as compared to other houses.

C. Outlier Deduction

We apply outlier deduction in the microgrid datasets. We apply all three algorithms discuss in the previous section. We have found similar results in all these algorithms. 10 homes have considered being the outlier out of 443 houses in this microgrid. The outlier houses id name as b159,b172,b44,b674,b688,b698,b706,b730,b803 and b879.b674 house having consuming to much power as compared to others. Their consumption affects the whole micro-grid. These outliers further can be used in clustering for best results.

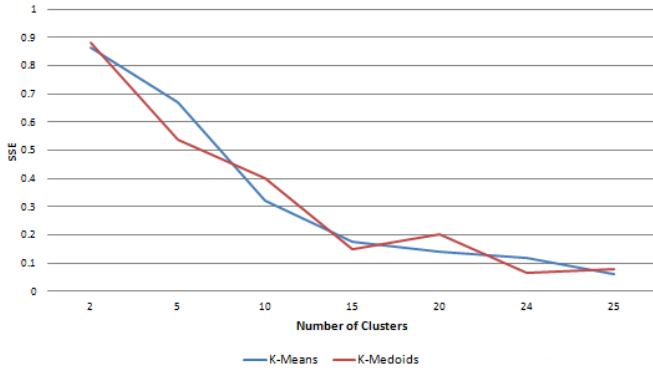


Fig. 4. SSE plot of House 1 REDD Datasets

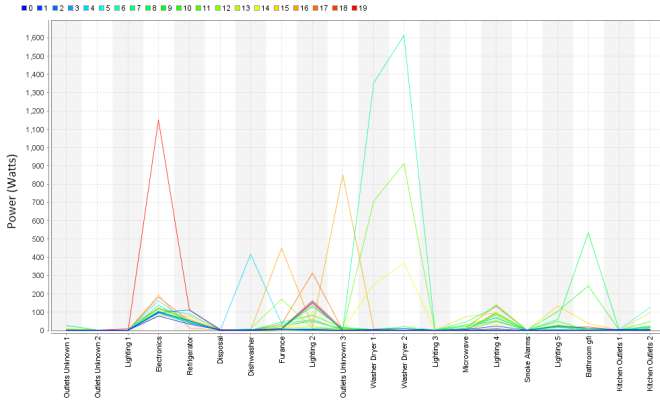


Fig. 5. Clusters of House 3 REDD Datasets

D. Clustering

We apply the clustering on both the datasets and achieved very interesting results. First, we apply the k-mean and k-medoids clustering algorithm on REDD datasets. Both the algorithm give almost similar clusters. SSE plot of both algorithm on House 1 load profile of REDD dataset in Figure 4. The plot shows that the SSE decrease when the number of cluster or k is almost equal to the number of devices use in the house. so we use k value is the number of devices monitor in each house. We also ignore the k-medoids on other datasets because they give similar result so we apply the k-means algorithm on all 6 houses of REDD datasets. The centroid table plot of 3rd house out of 6 houses shown in Figure 5. After applying clustering algorithm we see that there are peaks in clustering plot these peaks is high consumption of devices uses in these houses. Devices that having more consumption in all 6 houses are dishwasher, wash dryer, electric heater, air condition and furnace. So clustering help us in this datasets to find which home devices take more consumption as compared to other devices. This would be help in peak analysis and also peak reduction.

In microgrid datasets, we apply the k-means algorithm on whole datasets and datasets without outliers. These datasets having some outlier which affects overall results. Figure 6 shows the clustering plot whole dataset and we can see that one cluster having a huge peak as compared to other clusters. Figure 7 shows the clustering plot without outlier and you can

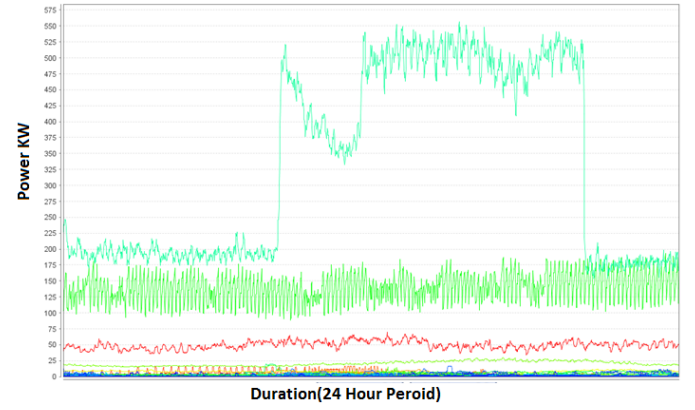


Fig. 6. Clusters of Microgrid Datasets

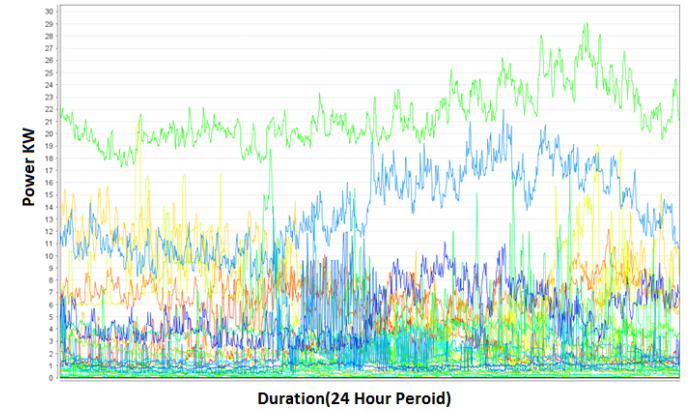


Fig. 7. Clusters of Microgrid Datasets without outliers

see significance difference between Figure 6 clusters. In this, all clusters are mostly close to each other. The results show that outlier effects the consumption of the whole microgrid. If we remove these outliers from microgrid, then the supplier can easily manage the production of the load.

E. Classification

We already discuss the classification help in the different area of research. We apply the classification on microgrid datasets by using a decision tree. This decision tree generates after applying the clustering on the datasets. Figure 8 shows the classification of microgrid datasets that having no outliers. We can see in the figure that decision tree classify the data into seven different branches. This classification helps for prediction or categorize the data.

F. Cluster Prediction

Clustering can also be help in prediction but in our datasets the clustering prediction result not much significance that can be shown in Table III. We apply the prediction on house 1 of REDD datasets per minutes and per seconds then apply on the microgrid. In House 1, our goal is to predict the load when different devices in on or off. To handle the on or off devices, we first convert the consumption value into binomial. Then we apply k-mean clustering and Map clustering on Label after that we have to find 63.28% accuracy on per

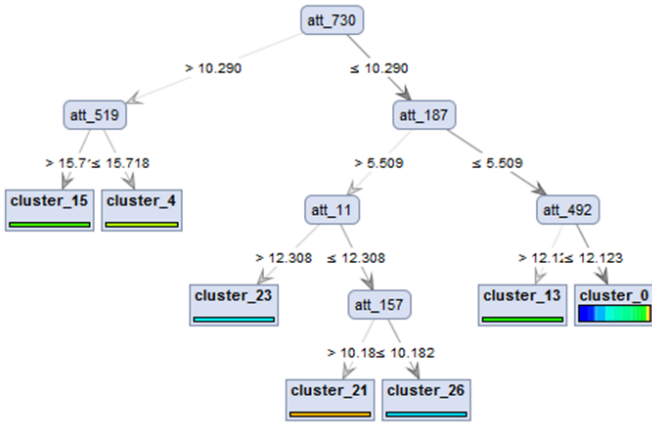


Fig. 8. Decision Tree of Microgrid datasets without outliers

seconds data and 62.31% accuracy on per minutes data. These results are not much good. We also apply the prediction on microgrid datasets. Before applying the prediction we discretize by binning the total consumption of microgrid that makes the ranges of consumption create labels for prediction. The accuracy of this data set is 50.35% which is also not good and less than REDD datasets.

TABLE III. CLUSTERING PREDICTION RESULTS

Dataset	Accuracy	Recall	Precision	F-Measure
House 1 per Seconds	63.28%	16.52%	11.49%	27.32%
House 1 per Minutes	62.31%	9.62%	5.59%	7.07%
Microgrid dataset	50.35%	46.35%	46.71%	46.53%

G. Frequent Pattern Mining

We apply the frequent pattern mining on the REDD datasets on all the houses. Our goal is to find most frequent devices use in these houses. First, we convert the data into the binomial format and then apply the FP-Growth algorithm on it. Table IV shows the most frequent devices used in all 6 houses and we can see that the Kitchen devices mostly use as compare to other devices like Refrigerator, Microwave, Stove and Lighting etc. The algorithm we use FP-Growth that already discuss in the previous section.

TABLE IV. FREQUENT DEVICES USE IN THE REDD DATASETS

House	Frequent Devices
1	Refrigerator, Microwave, Lighting 3, Kitchen Outlets 3, Kitchen Outlets 2, Kitchen Outlets 1, Bathroom gfi
2	Refrigerator, Microwave, Lighting, Kitchen Outlets 2
3	Refrigerator, Microwave, Furnace, Electronics
4	Stove, Miscellaneous, Lighting 2, Kitchen Outlets 2, Kitchen Outlets 1, Furnace
5	Subpanel 2, Subpanel 1, Refrigerator, Outlets Unknown 1, Microwave, Lighting 4, Lighting 3, Lighting 1, Electronics, Bathroom gfi
6	Washer Dryer, Stove, Refrigerator, Lighting, Kitchen Outlets 2, Kitchen Outlets 1, Electric Heat, Bathroom gfi, Air Conditioning 1

IV. CONCLUSION

In this paper, we discuss in detail major data mining techniques for understanding electricity load profile. After

applying these techniques on two different types of datasets we can easily understand the significance of data. Also analyze how we can improve the new power system by understanding large scale load profile data. All these techniques useful in the peak load reduction. The detail household level datasets understanding which household level devices take more load consumption as compared to other. That will also helpful in demand side management. In future, we can improve the results by applying these results in real time or more detail level datasets like seasonal and user socio-demographic information.

REFERENCES

- [1] Rakesh Agrawal, Tomasz Imieliński, and Arun Swami. Mining association rules between sets of items in large databases. In *ACM SIGMOD Record*, volume 22, pages 207–216. ACM, 1993.
- [2] Sean Barker, Aditya Mishra, David Irwin, Emmanuel Cecchet, Prashant Shenoy, and Jeannie Albrecht. Smart*: An open data set and tools for enabling research in sustainable homes. *SustKDD*, August, 2012.
- [3] Markus M Breunig, Hans-Peter Kriegel, Raymond T Ng, and Jörg Sander. Lof: identifying density-based local outliers. In *ACM sigmod record*, volume 29, pages 93–104. ACM, 2000.
- [4] Daswin De Silva, Xinghuo Yu, Daminda Alahakoon, and Grahame Holmes. A data mining framework for electricity consumption analysis from meter data. *Industrial Informatics, IEEE Transactions on*, 7(3):399–407, 2011.
- [5] Ian Dent, Uwe Aickelin, and Tom Rodden. The application of a data mining framework to energy usage profiling in domestic residences using uk data. *arXiv preprint arXiv:1307.1380*, 2013.
- [6] Dipl-Wi-Ing Christoph Flath, Dipl-Wi-Ing David Nicolay, Tobias Conte, PD Dr Clemens van Dinther, and Lilia Filipova-Neumann. Cluster analysis of smart metering data. *Business & Information Systems Engineering*, 4(1):31–39, 2012.
- [7] Jiawei Han, Jian Pei, and Yiwen Yin. Mining frequent patterns without candidate generation. In *ACM SIGMOD Record*, volume 29, pages 1–12. ACM, 2000.
- [8] Luis Hernández, Carlos Baladrón, Javier M Aguiar, Belén Carro, and Antonio Sánchez-Esguevillas. Classification and clustering of electricity demand patterns in industrial parks. *Energies*, 5(12):5215–5228, 2012.
- [9] AM Jorgensen and H Karimabadi. A survey of data mining techniques. In *AGU Fall Meeting Abstracts*, volume 1, page 02, 2005.
- [10] Leonard Kaufman and Peter Rousseeuw. *Clustering by means of medoids*. North-Holland, 1987.
- [11] J Zico Kolter and Matthew J Johnson. Redd: A public data set for energy disaggregation research. In *Workshop on Data Mining Applications in Sustainability (SIGKDD)*, San Diego, CA, volume 25, pages 59–62. Citeseer, 2011.
- [12] James MacQueen et al. Some methods for classification and analysis of multivariate observations. In *Proceedings of the fifth Berkeley symposium on mathematical statistics and probability*, volume 1, pages 281–297. Oakland, CA, USA., 1967.
- [13] Fintan McLoughlin, Aidan Duffy, and Michael Conlon. A clustering approach to domestic electricity load profile characterisation using smart metering data. *Applied energy*, 141:190–199, 2015.
- [14] Thair Nu Phyu. Survey of classification techniques in data mining. In *Proceedings of the International MultiConference of Engineers and Computer Scientists*, volume 1, pages 18–20, 2009.
- [15] B Pitt. Applications of data mining techniques to electric load profiling. *Dept. Elect. Electron. Eng*, pages 1–197, 2000.
- [16] Tri Kurniawan Wijaya, Tanuja Ganu, Dipanjan Chakraborty, Karl Aberer, and Deva P Seetharam. Consumer segmentation and knowledge extraction from smart meter and survey data. In *SDM*, pages 226–234. SIAM, 2014.