

Feature selection for the non-intrusive electrical appliances identification

Piotr Bilski¹, Wiesław Winiński²

¹*Institute of Radioelectronics and Multimedia Technologies, ul. Nowowiejska 15/19, 00-665
Warsaw, Poland, pbilski@ire.pw.edu.pl, +48601056969, pbilski77*

²*Institute of Radioelectronics and Multimedia Technologies, ul. Nowowiejska 15/19, 00-665
Warsaw, Poland, W.Winiński@ire.pw.edu.pl, +48 22 234-7341*

Abstract – The paper presents the application of selected feature selection methods to determine their optimal set in the task of non-intrusive appliances identification. The process detects devices operating in the specific location based on the analysis of voltage and current signals. The informativeness calculation methods and Fisher linear discriminative analysis are proposed to find significant features for determining, which appliances changed their state. The approach is verified on two sets of features. The scheme was tested in the laboratory conditions with multiple appliances working simultaneously, exploiting the decision tree as the appliance identifier. Experimental results show usefulness of the approach maximizing the identification accuracy with decreased complexity of computations.

Keywords – NIALM, artificial intelligence, feature selection

I. INTRODUCTION

The Non-Intrusive Appliance Load Monitoring (NIALM) is the hot topic in the green economy. The aim to suppress power consumption in households requires application of sophisticated approaches for the analysis of the aggregated current and voltage signals, recorded at the single location, near the energy meter. Multiple methods were proposed to identify appliances operating in the apartment. They differ by the form of presented knowledge, frequency range of analysed signals and features used to make decisions. Because the set of available appliances is vast, maintaining the high identification accuracy is difficult, justifying the application of various methods and their combinations. The prominent role here is played by the Artificial Intelligence (AI) approaches, able to autonomously process available data and make decisions with the minimum error. Their advantages include the generalization (proper responses for the new data, not

experienced before) and the ability to learn, i.e. extract knowledge from the present data. Their accuracy depends on the processed features, which should facilitate distinguishing between various appliances.

The paper presents the approach to analyse and select the most important features from the appliance identification standpoint, extracted from the current and voltage signals. Three methods were selected to minimize the number of characteristic values to be calculated before the decision is made. These include two informativeness measures and the Fisher discriminative analysis criterion. Their application simplifies the data acquisition process, while maintaining as high identification accuracy, as possible.

The structure of the paper is as follows. In section II the principles and aims of the applied NIALM system are introduced. Section III describes the feature selection methods applied for the task. In section IV results of the methods' verification are presented. Section V contains conclusions and future prospects.

II. ARCHITECTURE OF THE NIALM SYSTEM

The proposed structure of the NIALM system used in the research is in Fig. 1. It consists of two parts, representing the hardware and software modules. The former is the data acquisition (DAQ) node and the computer running the appliance identification software. The program is run online, analysing the voltage and current samples provided by the DAQ node. Its tasks include:

- detecting the change in the current level, which is the premise that the configuration of operating appliances was changed.
- calculating features required by the classifier
- identifying the appliance that changed its state.

The system operates online, requiring the fast DAQ card to collect samples and the processing power-efficient software. It is then important to decrease the number of calculated features. The current and voltage signals are sinusoidal, with the frequency of 50Hz. The selected

sampling frequency is the key parameter influencing the set of features that can be computed. In the literature, three ranges are discussed: single Hertz [1], tens [2] and hundreds (or even thousands) [3] of kHz. In the presented case the second option was selected, where $f_s=2kHz$ is enough to monitor the first twenty harmonic components.

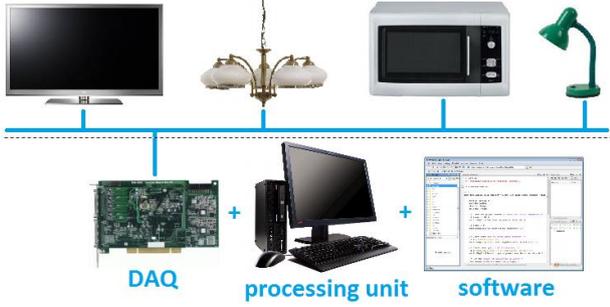


Fig. 1. Structure of the NIALM system

The software operates on the principle that each appliance requires the particular amount of current during the operation. Therefore turning it on results in the increase of the consumed current, while turning off leads to its decrease. Except the two-state devices (such as the bulb or the kettle), there are also multi-state appliances (washing machine or the microwave), running multiple programs. In their case multiple current levels should be considered. Such analysis is more complex [4] and requires individual approach. In the particular state the current consumption is constant. The analysis of the acquired waveforms focuses on their envelopes (created by the maximum values of each samples' vector with the predefined length). After determining the change in the current level (further called the "event"), the identification is commenced. The first parameter of the event detection module is the minimum current change, suggesting the change in the appliances' configuration (some device was switched on or off). It was set to $200mA$ as the compromise between the ability to detect power saving devices and the number of detectable events. The second parameter is the frequency of checking for the occurrence of the event. It was assumed that the devices' configuration may change twice a second at best.

The appliance identification procedure is run after detecting the event. The AI-based classification method (such as rule-based approaches [5] or artificial neural networks [6]) obtains the vector of features to determine which device changed its state regarding the previously calculated configuration. The classification efficiency depends on the event detection procedure. If one event is missed, the appliance identification will make a mistake (by ignoring one appliance). On the other hand, false events are possible, where no actual change took place (if the detection threshold is too low). Knowledge used for the identification is obtained from training data sets, which must be provided prior to the online system application. The preparation of the set consists in collecting multiple feature vectors for each appliance, made after monitoring

it for the predefined amount of time with all other appliances turned off. The on-line identification works with multiple appliances turned on, working on additive features. After detecting the event, features are extracted from the actual current waveform and subtracted from the vector calculated after the previous event. This way the information about the specific appliance is extracted.

III. FEATURE SELECTION METHODS

Depending on the sampling frequency, various sets of features f_j are considered. Time and frequency domain parameters are selected, including the power characteristics and harmonic components [2]. The open question is their significance for the device identification. When no a priori knowledge about their importance exists, the safe way is to generate their maximum number and provide to the AI classifier during the machine learning (ML). In the process, only the most significant attributes will be selected. If only important ones are present in the data set S , the ML is faster and calculations in the online mode are limited to only the subset of features. The dimensionality reduction is presented in Fig. 2, where two algorithms are exploited: one for reducing the set S , the second for verifying the reduction efficiency. The classification accuracy is expected to be the same for both cases, but with shorter training time for the second set.

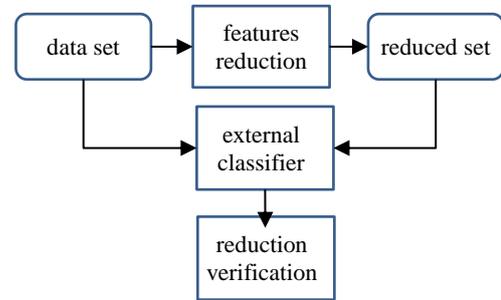


Fig. 2. Operation of the data set reduction

This section introduces applied feature selection methods. Vectors (examples e_i) of collected features from which the optimal set is selected, are presented. The decision tree (DT) classifier, working on the original and minimized sets is discussed. The processed data have the form of the table (1) with n examples, m features and k integer categories of appliances $c=\{1, \dots, k\}$ (supplemented by vectors calculated for no appliance operates, required to identify the "no change" event, indicated by 0 value).

$$S = \begin{bmatrix} e_1 \\ \vdots \\ e_n \end{bmatrix} = \begin{bmatrix} f_1 & c_1 \\ \vdots & \vdots \\ f_n & c_k \end{bmatrix} = \begin{bmatrix} f_{11} & \dots & f_{1m} & c_1 \\ \vdots & \ddots & \vdots & \vdots \\ f_{n1} & \dots & f_{nm} & c_k \end{bmatrix}, \quad (1)$$

A. Feature selection strategy

The approaches used to analyse features are well known,

although they are usually used in the situation of only two categories considered. The presented case refers to the multi-category situation (each appliance is represented by the separate category) these approaches had to be modified. After the calculation, the obtained set of measure values $\mathbf{q}=\{q_1, \dots, q_m\}$ for the specific appliance is ordered into a sequence $o(\mathbf{q})=\{o_1(\mathbf{q}), \dots, o_m(\mathbf{q})\}$, from the most to the least important. The measures are computed for every appliance separately. Therefore for k appliances there are k sets of ordered features, in each case with different sequence. After calculating the measures, one resulting set has to be constructed, containing features important for all appliances. This requires determining how many features of the ordered sequences should be considered for processing. Three strategies (further referred to as s_a , s_b and s_c respectively) are proposed here:

- a) Selection of the top number of features that are separated from the others by the greatest difference between the neighbouring values. For the ordered set of significance measures, the differences $\Delta o(\mathbf{q})$ (2) are calculated and the largest one is the cut-off point. This is the reasonable choice during the selection step. Unfortunately, in some cases it leads to large sets of features and low reduction rate.

$$\{\Delta o(\mathbf{q})\} = \{o_1(\mathbf{q}) - o_2(\mathbf{q}), \dots, o_{m-1}(\mathbf{q}) - o_m(\mathbf{q})\} \quad (2)$$

- b) Selection of the predefined number t of top features, eliminating too numerous groups. In the presented research the threshold depends on the number of features and is set at $\lceil 0.09 \cdot m \rceil$.
- c) Selection of the number of top features between t and h to avoid both too long and too short groups for the particular appliances. In the presented research the fraction of selected features was set between $\lceil 0.05 \cdot m \rceil$ and $\lceil 0.07 \cdot m \rceil$. This allows for creating relatively small (but not too small) groups.

B. Selection algorithms

The first method is the Fisher linear discriminative analysis, proposed in [7] for the binary case and expanded to cover multiple categories in [8]. The idea was to calculate the significance measure q_{jk} for each pair of features (resulting in $(j-1)!$ vectors). The approach proposed here allows for calculating only j measure vectors (3), which can be further ordered sequentially.

$$q_F(j) = \frac{\text{abs}(\mu_{j,l} - \mu_{\sim j,l})}{\text{abs}(\sigma_{j,l} - \sigma_{\sim j,l})}, \quad (3)$$

Here $\mu_{j,l}$ is the mean value of the j -th feature for examples representing the appliance l , $\mu_{\sim j,l}$ is similarly the mean value of the j -th feature for examples belonging to all other categories. The values $\sigma_{j,l}$ and $\sigma_{\sim j,l}$ are standard deviations of features interpreted analogously. The *abs* operator stands for the absolute value. This way, instead of

calculating measures for all the pairs, they are computed between the actual appliance category and all other.

The second method is the informativeness measure, based on the correlation between the particular features in examples belonging to different categories. The formula (4) expresses the measure for the single feature, where $||$ is the number of examples in the set S that belong (or not) to the selected category j .

$$q_{i1}(j) = \frac{\text{abs}(\mu_{j,i})}{\sqrt{\frac{\sum_{\sim j} \frac{1}{|c(e_i) \neq j|} \sum_{\sim j} b_j \cdot |c(e_i) = j| - 1}}{m-2}}, \quad (4)$$

The b_j coefficient is calculated as in (5), where $\bar{\mu}_i$ is the difference between the mean value of the j -th feature calculated for examples belonging to the j -th category and all others.

$$b_j = \frac{1}{|c(e_i) = j| - 1} \cdot \sum_j (f_{ji} - \bar{\mu}_i)^2, \quad (5)$$

This way the significance of the set of features to distinguish between the j -th category from all others ($\sim j$) can be calculated based on the scatter of the features' values around their means.

Alternatively, the informativeness can be calculated as the scalar product for examples belonging to different categories, treated as the vectors in the m -dimensional space. Feature vectors are gathered into two groups: belonging to the particular category or to all others, with the same aim as before.

$$q_{i2}(j) = \frac{(f_i \cdot c(e_i))^2}{\|f_i\|^2 \cdot \|c(e_i)\|^2} = \frac{(f_i \cdot c(e_i))^2}{\sqrt{\sum_{i=1}^n f_i^2} \cdot \sqrt{\sum_{i=1}^n c_i}}, \quad (6)$$

Although the presented measures are used as the standard tool in the data analysis, their comparison is required, because they all produce different ordered sequences of features. It is assumed that relations between the features and categories are linear, which does not have to be true. The presented approaches were implemented to verify their usefulness to the NIALM.

C. Collected features

In the medium range of frequencies (up to tens of kHz) various features may be extracted from the current and voltage waveforms. Two sets of features were prepared for the processing. Based on the preliminary analysis it was assumed that the first one is required to perform the accurate identification, while the second plays the auxiliary role in the process. The set S_j contains the following 67 features (the brackets contain their numbers):

- Mean current value (1)
- DC value (2)
- Mean power (3)

- The first 16 harmonics of active power (4-19)
- The first 16 harmonics of reactive power (20-35)
- The first 16 harmonics of the real current components (36-51)
- The first 16 harmonics of the imaginary current components (52-67)

The set S_2 was introduced to gain independence of the measurement process from the voltage disruption introduced by external power networks. It contains sixteen components (numbered 1 to 32 or 68 to 99, depending on whether S_1 is considered separately or as the part of $S=S_1+S_2$):

- The first 16 harmonics of conductances (1-16 or 68-83)
- The first 16 harmonics of susceptances (17-32 or 84-99)

The experiments consisted in applying the identification algorithm (DT) processing three sets of features (i.e. S_1 , S_2 and S_1+S_2) without and with the feature reduction. This way it is possible to determine, which features are the most important for the appliances identification and what algorithm is the most suitable for this purpose.

D. Verification algorithm

The DT classifier is the simple, memory efficient method, used in multiple applications, including NIALM [9]. The proposed version of the inductive tree generation was applied in [10], leading to various variants of the knowledge structure, which depend on the method of selecting the tests for the node. The tree (Fig. 3) consists of the nodes used to direct the analysed set of features (testing example) to one of leaves (blue terminal nodes) containing the appliance identifier that changed its state during the event (which is presented as the decision made by the DT). The nodes contain the test (the feature and its value). Based on the comparison between the example's feature and its threshold value θ_i in the node, the example is directed to one of two subtrees, until reaching one leaf.

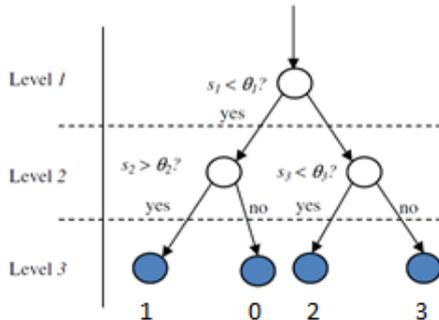


Fig. 3. Structure of the decision tree

The advantage of the DT is the legible form of the stored knowledge. The analysis of its structure leads to the information about which features were selected during the learning process from the training set S . The applied

inductive DT generation (based on the C4.5 algorithm [11]) consists in adding nodes to the tree with the tests created to separate different appliance identifiers as quickly as possible. This leads to the simplest knowledge structure. The candidates to the test are calculated during each algorithm iteration as the middle values between every pair of the ordered j -th feature values from S .

$$a_{ij} = \Delta f_{ij} = f_{ij} - f_{i-1,j}, \quad (7)$$

The best ones are automatically selected according to the entropy criterion, but the problem arises when there are multiple features' values with identical, minimum entropy. The following actions were applied in such a situation (italics in brackets indicate the abbreviation of the strategy, used further). The value of the feature is selected so that:

- the distance from the neighbouring values (7) is maximal (*maxdist*),
- the distance from the neighbouring values (7) is minimal (*mindist*),
- its number of occurrences in the tree is maximal (*maxocc*),
- its number of occurrences in the tree is minimal (*minocc*),
- the choice is random (*random*).

IV. EXPERIMENTAL SCHEME

The presented approach was verified for the laboratory test stand, containing six two-state appliances (turned on or off). They were (the appliance identifier is in brackets): power saving bulb (1), dryer (2), vacuum cleaner (3), mixer (4), juicer (5) and kettle (6). Multiple sequences were created, leading to the series as in Fig. 4. The RMS current value visible on the upper part is processed to detect the events, while lower part illustrates the actual sequence of events during the half-hour monitoring, during which both current and voltage waveforms were collected.

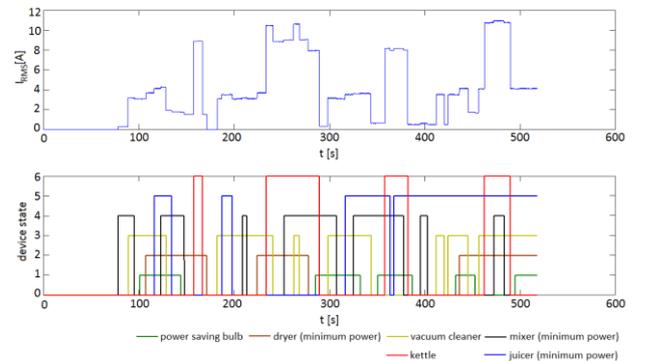


Fig. 4. Sequence of the analyzed appliances operation

The first experiment involved training the DT on three data sets to compare their usefulness for the knowledge extraction. The training was performed on the sets containing one hundred feature vectors calculated for each

appliance operating with other devices switched off. After including the measurements for all appliances switched off, the set S contained 700 examples. The testing was performed on the sequence from Fig. 4. In Tab. 1 results of the identification for three sets are presented. The accuracy percentage is the relative number of correctly identified categories (including the “no change” event).

Table 1. Identification results [%] depending on the processed feature sets.

Test selection	S_1	S_2	S_1+S_2
<i>mindist</i>	42.28	91.42	45.14
<i>minocc</i>	44.57	36.57	44.57
<i>maxocc</i>	47.42	90.85	47.42
<i>maxdist</i>	92.57	46.85	92.00
<i>random</i>	52.57	21.71	56.00

The set S_1 contains the most useful features, which is indicated by the maximal identification accuracy (marked by the bold font). As could be predicted, the DT with the *maxdist* strategy was the most useful (intuitively, it is the method ensuring generalization in the widest sense). The auxiliary set S_2 leads to a little worse outcomes, although in the case of external distortions it should be preferred. The combined set S_1+S_2 does not increase the identification accuracy, therefore there is no point in using the whole spectrum of features. The random selection (which results are mean values of 10 repeated DT generations) of the node test leads to poor performance and should be rather implemented in the random forest (RF).

Table 2. Identification outcomes for the optimal DT processing the set S_1 .

Event No.	Predicted	Actual
1	0	0
31	6	6
33	3	3
45	4	4+1
51	4	0
75	6	6+1
99	0	1
164	1	1

The accuracy was calculated for 175 events (see Tab. 2), most of which were false alarms, correctly classified in majority of cases (like event No. 1). The widespread problem was distinguishing between the “no change” and the power saving bulb (event No. 99) or the mixer (event No. 51). It is difficult to avoid the former, especially when the bulb changes its state with the large current consumption in the background (events No. 45 and 75). In such a case, the DT detects the change of the more prominent appliance, although the bulb is missed by the event detection module. The incorrect identification of low power devices is not a problem, as such appliances are barely visible in the overall power consumption pattern. Missing more significant devices leads to the

inconsistencies in the detected configuration and the overall power consumption (measured independently). The latter must be used periodically to correct such errors.

Application of the feature selection methods allowed to obtain reduced subsets of all three sets, i.e. S_1 , S_2 , S_1+S_2 . The ordering results for the Fisher criterion performed on the set S_1 for the power-saving bulb are in Fig. 5. The feature indexes are defined as in section III.C. In most cases, the first two or three features are considered much more important than all others. This is not the case for the “no change” state, leaving the large number of attributes. The selection strategy s_a leads to the overall set of 45 features using q_f , 67 for q_{i1} and 47 using q_{i2} . The classification results for such subsets are the same as in Tab. 2, as all important features are in the reduced sets. This amount is not acceptable due to the requirement of decreasing the computational power during the features calculation. Therefore the selection strategies s_b and s_c were used for the set S_1 , leading to the significant decrease of the number of features. For the limitations introduced in section III.A, the strategy s_b imposed the maximum number of features for the single appliance to 6. For the strategy s_c the range of remaining features was between 3 to 5. Results of the DT operation in both situations are in Tab. 3, sets of remaining features are in Tab. 4.

Table 3. Identification results [%] for the optimal DT processing the reduced set S_1 .

Test selection	s_b			s_c		
	q_f	q_{i1}	q_{i2}	q_f	q_{i1}	q_{i2}
<i>mindist</i>	42.28	89.71	45.14	42.28	89.71	45.14
<i>minocc</i>	44.57	39.92	44.57	44.57	39.92	44.57
<i>maxocc</i>	47.42	89.85	47.42	47.42	89.85	47.42
<i>maxdist</i>	92.57	46.85	92.57	92.57	46.85	92.57
<i>random</i>	54.85	42.85	56.00	54.85	42.85	56.00

Table 4. Reduced feature sets' structure for various strategies.

Strategy	Feature selection method	Remaining features
s_b	q_f	2, 3 , 4 , 11, 13, 14, 21, 28, 29, 36, 37, 43, 46, 53, 60, 61
	q_{i1}	8, 12, 16, 23, 28, 29, 32, 33, 35, 40, 44, 47, 48, 49, 51, 55, 60, 61, 63, 64, 65, 66, 67
	q_{i2}	1, 2, 3 , 4 , 11, 12, 14, 21, 22, 26, 28, 36, 43, 44, 46, 53, 54, 58, 60
s_c	q_f	1, 2, 3 , 4 , 5, 9, 11, 12, 13, 14, 21, 28, 29, 36, 37, 43, 44, 46, 53, 60, 61, 62
	q_{i1}	11, 12, 14, 16, 19, 28, 29, 32, 33, 34, 35, 43, 44, 47, 48, 51, 60, 61, 62, 63, 64, 65, 66, 67
	q_{i2}	1, 2, 3 , 4 , 5, 11, 14, 21, 22, 26, 28, 36, 37, 40, 43, 45, 46, 53, 54, 58, 60

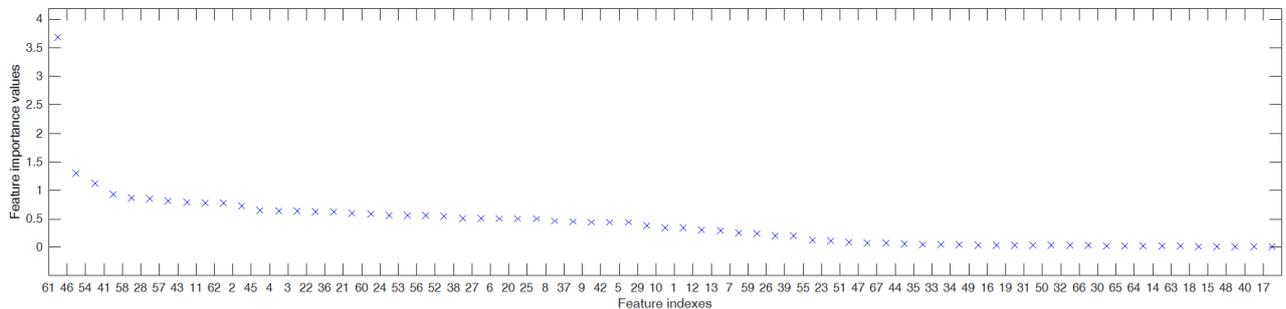


Fig. 5. Results of the features' importance ordering for the power-saving bulb (Fisher criterion)

The analysis of the optimal DT structure shows that it uses the mean power (3), the first harmonics of the active (4) and reactive (20) power. The first two features are used in the high level of the tree and are the most important. As long as they are present in the reduced set (in Tab. 4 indicated by the bold font), the identification accuracy remains unchanged (the third feature may be replaced by some other). The Fisher and the second informativeness criteria give similar results and can be used interchangeably, while the first informativeness criterion does not reflect the importance of subsequent features correctly. After reducing features to about 25% of the original set, the identification quality is maintained.

V. CONCLUSIONS

The proposed methodology allows for significantly decreasing the number of analysed features, required to correctly identify changes in the configuration of appliances operating in the apartment. The algorithms used lead to the computationally effective approach, both in the online and offline mode. However, there is the need to further suppress the number of selected features. Also, the verification method can be replaced by another algorithm (such as the RF or support vector machines).

The future research should cover introduction of other significance measures, considering the nonlinear relations between the features (such as the Kolmogorov-Smirnov test [12]). The more complex appliances, such as the finite state machines (washing machines) should be analyzed.

VI. ACKNOWLEDGMENT

This work has been accomplished within the project PBS2/A4/0/2013 "The Non-invasive System for Monitoring and Analysis of Electricity Consumption in the Area of the End-user" financially supported by the National Centre for Research and Development.

REFERENCES

[1] G.W. Hart, „Nonintrusive Appliance Load Monitoring,” *Proceedings of the IEEE*, Vol. 80, No. 12, pp. 1870-1891, 1992.

[2] L. Jiang, S. Luo and J. Li, "Automatic power load event detection and appliance classification based on power harmonic features in nonintrusive appliance load monitoring," *Proc. 8th IEEE Conf. Industrial Electron. and App.*, pp. 1083-1088, 2013.

[3] S. Gupta, M. S. Reynolds and S. N. Patel, "ElectriSense: Single-Point Sensing Using EMI for Electrical Event Detection and Classification in the Home," *Proc. Ubicomp '10 Proceeding0073 of the 12th ACM int. conf. ubiquitous computing*, pp. 139-148, 2010.

[4] M. Zeifman, K. Roth, "Nonintrusive Appliance Load Monitoring: Review and Outlook," *IEEE Transactions on Consumer Electronics*, Vol. 57, No. 1, pp. 76-84, 2011.

[5] B. Qin, Y. Xia, S. Prabhakar, and Y. Tu, "A Rule-Based Classification Algorithm for Uncertain Data," *IEEE 25th International Conference on Data Engineering*, March 29 2009-April 2 2009, Shanghai, China, pp. 1633-1640.

[6] T.H. Oong, N.A.M. Isa, "Adaptive Evolutionary Artificial Neural Networks for Pattern Classification," *IEEE Transactions on Neural Networks*, Vol. 22, No. 11, Nov. 2011, pp. 1823-1836.

[7] S. Mika, G. Ratsch, J. Weston, B. Scholkopf, and K. R. Mullers "Fisher discriminant analysis with kernels," *IEEE Conference on Neural Networks for Signal Processing IX*, pp. 41-48, 1999, DOI:10.1109/NNSP.1999.788121.

[8] Y. Xu, F. Li, and T. Hu, "A Method of Kernel Fisher Discriminant for Multi-class Classification," *6th World Congress on Intelligent Control and Automation*, Vol. 2, pp. 9954 - 9957, 2006.

[9] M. Berges, E. Goldman, H.S. Matthews, L. Soibelman, "Learning Systems for Electric Consumption of Buildings," *Proceedings ASCE International Workshop on Computing in Civil Engineering*, 2009.

[10] P. Bilski, P. Mazurek, and J. Wagner, W. Winiiecki, "Application of decision trees to the fall detection of elderly people using depth-based sensors," *Proc. XXI IMEKO World Congress*, 2015.

[11] J. R. Quinlan. Improved use of continuous attributes in c4.5. *Journal of Artificial Intelligence Research*, 4:77-90, 1996.

[12] R. Marsalek and K. Povalac, "Kolmogorov-Smirnov Test for Spectrum Sensing: From the Statistical Test to Energy Detection," *Proceedings IEEE Workshop on Signal Processing Systems*, pp. 97-102, 2012. DOI: 10.1109/SiPS.2012.58.