

Control List Generation to Reveal Fraudulent Customers for Electricity Provider Companies

I. Harmati, G. Kovács, B. Kiss, G. Vámos and J. Fodor

Abstract— This paper presents a procedure that generates control lists for electricity provider companies to reveal fraudulent customers with the highest possible efficiency. Based on the database of Customer Registration System, essential features are selected and other secondary (derived) properties are introduced to build up an abstract space for the algorithm of control list generation. Control list generation algorithm uses training data from the Customer Registration System to set up a fuzzy system with subtractive clustering. The algorithm has been validated at Hungarian Electricity Provider ELMŰ-ÉMÁSZ Company Group.

I. INTRODUCTION

FRAUDULENT customers induce significant costs for service providers. It is especially true for electricity provider companies since electricity is the most common form of energy to run households and factories. Fraudulence can be originated from different kinds of behaviors, income level, social environment and even culture. In some regions or countries, the rate of fraudulent customers is high enough that the loss of electricity providers exceeds the cost of developing intelligent algorithms revealing the fraudulence. In this case, electricity providers are motivated to break down the number of fraudulence by controlling the customers on the top of control list (evaluated as highest chance for illegal use of service).

It is not easy to reveal customers who use electricity illegally. There are no strict rules or phenomena that guarantee the discovery of all the fraudulent customers. Service providers rely mostly on their experienced employees to create a control list. Since artificial intelligence methods are in close relation with human logic, a straightforward idea is to set up an inference system which works similarly to the way of human thinking [4], [5], [8]. Several artificial intelligence algorithms are able not only to learn human reasoning but also to discover relations hidden for human operators [10], [11].

This work was supported in part by ELMŰ-ÉMÁSZ Company-Group, the Hungarian National Scientific Research Foundation grant OTKA K71762. Also, it is connected to the scientific program of the "Development of quality-oriented and harmonized R+D+I strategy and functional model at BME" project, supported by the New Hungary Development Plan (Project ID: TÁMOP-4.2.1/B-09/1/KMR-2010-0002).

I. Harmati, G. Kovács, B. Kiss and G. Vámos are with Dept. of Control Engineering and Information Technology, Budapest University of Technology and Economics, Magyar Tudósok krt 2, Budapest, Hungary, H-1117 (e-mail: {gkovacs,vamos,bkiss,harmati}@iit.bme.hu).

J. Fodor is Head of Metering and Meter Checking Department - ELMŰ Network Distribution Ltd., Váci út 72-74, Budapest, Hungary, H-1132

The first step to apply machine learning is data analysis and data matching. The database of Customer Registration Systems of service providers is usually not applicable directly to any soft computing algorithm. As a result, the comprehensive database should be filtered, reorganized and reduced to build up a consistent and clarified database. To do so, one can use the methodologies of knowledge engineering process and knowledge modeling, and set up the work-flow and use-case diagrams for the Customer Registration System [1], [2], [3], [9]. A consistent database allows defining key variables, features and secondary (derived) properties that may play important role in reasoning. The experience of employees in charge of customer control at the service provider can be integrated into the algorithm at this point. Variables defined in this way compose the inputs of the inference system.

Several options exist to set up an artificial intelligence based inference system to identify fraudulent the customers. The method presented in this paper focuses on subtractive clustering method. The idea behind applying clustering method is that clustering methods are able to recognize connections between input and output variables [6], [7], [13]. They work reliably in practice and provide good general result even for multi input multi output systems. Our method adopts the subtractive clustering technique which is considered one of the most practical approaches. Subtractive clustering method delivers a Sugeno-type fuzzy system. Fuzzy systems have the advantage that their rules are close to the way of human thinking. In addition, the rules of fuzzy system with Sugeno-type can be tuned on-line as Adaptive-Neuro-Fuzzy Inference System (ANFIS) [12] when new data arrives into the system. The initial rules of Sugeno-type fuzzy system were defined on a subset of validated data. Validated data means in this case that customers in this dataset were checked for fraudulence therefore they are suitable for training.

The remaining part of the paper is organized as follows. Section II introduces Sugeno-type fuzzy systems and provides a brief overview about the algorithm of subtractive clustering used in our study. Section III summarizes the technical procedure carried out on the raw database in order to achieve a clarified consistent data set for further processing. Section IV defines the inputs, outputs and parameters of the fuzzy system and the subtractive clustering algorithm. Section V evaluates the efficiency of the method in comparison the conventional approach. Evaluation is carried out on real-life data from the Customer Registration

Systems. Section VI draws some conclusions and outlines possible research directions in the future.

II. THEORETICAL BACKGROUND

In this section, a brief overview is provided about the algorithms applied to reveal fraudulent customers. Namely, Sugeno-type fuzzy systems, subtractive clustering and its possible extension will be discussed in the sequel [5].

A. Sugeno-type fuzzy systems

The TSK- (Takagi-Sugeno-Kong) type or most often known as Sugeno-type fuzzy system unifies the principles of fuzzy and classical systems in such a manner which allows specifying deterministic consequence at the place of the consequence part of various relations [5]. This allows using for example different deterministic algorithms in various operating points. Of course in this case the operating point is not crisp, but fuzzy, thus various operating points can exist simultaneously with different possibilities. The knowledge base (rule base) of a TSK-type fuzzy system consists of relations in the following form:

$$R_i : \text{if } x_i \text{ is } X_1^i \text{ and } \dots \text{ and } x_N \text{ is } X_N^i \\ \text{then } y = f_i(x_1, \dots, x_N), \quad i = 1, \dots, n.$$

If $f_i = c_i$ is a constant function then the fuzzy system is a zero order Sugeno system. If

$$f_i = c_{i1}x_1 + \dots + c_{iN}x_N + c_{i0}$$

is a linear function then the fuzzy system is called first order Sugeno system. Sugeno-type fuzzy inference and defuzzification algorithm consist of two steps:

1) Inference:

$x_1^*, x_2^*, \dots, x_N^*$ are crisp inputs

$$\forall R_i \mapsto (\tau_i, y_i^*)$$

$$\tau_i = \mu_{X_1^i}(x_1^*) \wedge \dots \wedge \mu_{X_N^i}(x_N^*)$$

$$y_i^* = f_i(x_1^*, \dots, x_N^*)$$

2) Defuzzification:

$$y_{TSK}^* = \frac{\sum_i \tau_i y_i^*}{\sum_i \tau_i}$$

The main benefit of Sugeno-type fuzzy system appears in deterministic function of the consequence part of the rules. A deterministic function (mostly zero or first order function) makes parameter tuning easier. As a result, Sugeno-type fuzzy system is prepared to on-line adaptation.

B. Subtractive clustering algorithm

Subtractive clustering algorithm [5] is a possible method to approximate an unknown function $y = f(x)$. In order to simplify the notation, we consider simple systems with only one input and one output. The method can be easily generalized (and actually used in the implementation) for multivariable systems. The rule base has the form

$$R_i : \text{if } x \text{ is } A_i \text{ then } y \text{ is } B_i, \quad i = 1, \dots, m.$$

The following questions arise:

- What is m , the number of relations?
- What are the centers of the fuzzy sets A_i and B_i ?
- What are the membership functions $\mu_{A_i}(x)$ and $\mu_{B_i}(y)$?

We might assume that the teaching data is given in the form of sample points (x_i, y_i) , $i = 1, \dots, N$. We have two options to choose the centers of the fuzzy sets: either to choose the centers from the data points or on a raster (grid line). The first choice gives preference to the prescribed data, while the second tries to increase the interpolation property. Second option in our algorithm, and it is assumed, that the raster centers $N_{ij} = (x_i^*, y_j^*)$ are chosen in the intersection of grid lines in a rectangle (in the general case in a hypercube) containing the teaching data. In order to answer the questions we have performed the following algorithm.

Subtractive clustering algorithm:

1. Quantization of x and y to the allowed raster centers, that is the choice of the raster $N_{ij} = (x_i^*, y_j^*)$ as intersections of gridlines.
2. Approximation of the density of the sample points. A potential (mountain) function $M : N_{ij} \rightarrow R^1$ will be determined for this purpose by the following rule:

$$d(N_{ij}, (x_k, y_k)) = (x_k - x_i^*)^2 + (y_k - y_j^*)^2,$$

$$M(N_{ij}) = \sum_{k=1}^N \exp\{-\alpha \cdot d(N_{ij}, (x_k, y_k))\}.$$

3. Initialization: $m := 1$, $M_1 := M$, and choice of α, β and the stop condition $\delta > 0$ ($M_{m+1}^* < \delta$).
4. Cycle:

$$i) \quad M_m^* := \max_{ij} M_m(N_{ij}),$$

$$N_m^* := \arg \max_{x_i^*, y_j^*} M_m(N_{ij}) = (\bar{x}_m^*, \bar{y}_m^*)$$

- ii) Potential function supervision (subtracting):

$$M_{m+1}(N_{ij}) := M_m(N_{ij}) - M_m^* \exp\{-\beta \cdot d(N_m^*, N_{ij})\}$$

iii) Jump to step 4 if $M_{m+1}^* \geq \delta$, otherwise stop.

The cluster centers are chosen as $N_i^* = (\bar{x}_i^*, \bar{y}_i^*)$, $i = 1, \dots, m$, and for every cluster center $N_i^* = (x_i^*, y_i^*)$ a relation is assumed as

$$R_i: \text{if } x \text{ is near } x_i^* \text{ then } y \text{ is near } y_i^*, \quad i = 1, \dots, m.$$

Let the centers of the fuzzy sets A_i and B_i be $\bar{x}_i := \bar{x}_i^*$ and $\hat{y}_i := \bar{y}_i^*$, respectively. We can choose Gaussian membership functions for A_i (characterizing 'near'):

$$\mu_{A_i}(x) = \exp\left\{-\frac{1}{2}\left(\frac{x - \bar{x}_i}{\sigma_i}\right)^2\right\}.$$

A usual choice for σ_i is the following:

$$\frac{1}{2\sigma_i^2} = \beta \Leftrightarrow \sigma_i := 1/\sqrt{2\beta}.$$

Considering only Sugeno-type fuzzy systems we can assume that B_i is a fuzzy singleton. Observe that we have initialized a zero order Sugeno system. The firing weights are $\tau_i := \mu_{A_i}(x)$ and after defuzzification the prescribed data are approximated by the function

$$\hat{y} = \hat{f}(x) = \frac{\sum_{i=1}^m \tau_i \hat{y}_i}{\sum_{i=1}^m \tau_i} = \sum_{i=1}^m \tau_i^* \hat{y}_i, \quad \tau_i^* = \frac{\tau_i}{\sum_{k=1}^m \tau_k}$$

It is easy to observe that subtractive clustering algorithm creates fuzzy rules in an automatic way, which means that the number of rules and parameters are not predefined.

C. Possible improvement of rule base

The parameters \bar{x}_i, σ_i and \hat{y}_i can be further tuned by optimum seeking methods based on the available training data set. The partial error and the total error are

$$E_k = \frac{1}{2} [y_k - \hat{y}(x_k)]^2, \quad E_{total} = \sum_{k=1}^N E_k$$

respectively. If p is a parameter to be tuned, then

$$\frac{\partial E_k}{\partial p} = -[y_k - \hat{y}(x_k)] \frac{\partial \hat{y}(x_k)}{\partial p},$$

$$\frac{\partial E_{total}}{\partial p} = \sum_{k=1}^N \frac{\partial E_k}{\partial p},$$

hence it is enough to deal with the computation of $\partial \hat{y}(x_k) / \partial p$. Let us consider the different parameters. The partial derivatives needed for the gradient vector can be determined in the following steps:

$$\frac{\partial \hat{y}(x_k)}{\partial \hat{y}_i} = \tau_i^*,$$

$$\frac{\partial \hat{y}(x_k)}{\partial \bar{x}_i} = \hat{y}_i \frac{\partial \tau_i}{\partial \bar{x}_i} \frac{1}{\sum_k \tau_k} + \sum_i \hat{y}_i \tau_i (-1) \frac{1}{(\sum_k \tau_k)^2} \frac{\partial \tau_i}{\partial \bar{x}_i},$$

$$\frac{\partial \tau_i}{\partial \bar{x}_i} = \exp\left[-\frac{1}{2}\left(\frac{x_k - \bar{x}_i}{\sigma_i}\right)^2\right] \cdot \left(-\frac{1}{2}\right) \frac{1}{\sigma_i^2} 2(x_k - \bar{x}_i)(-1) = \tau_i \frac{x_k - \bar{x}_i}{\sigma_i^2},$$

$$\frac{\partial \hat{y}(x_k)}{\partial \bar{x}_i} = [\hat{y}_i - \hat{y}(x_k)] \tau_i^* \frac{x_k - \bar{x}_i}{\sigma_i^2},$$

$$\frac{\partial \hat{y}(x_k)}{\partial \sigma_i} = [\hat{y}_i - \hat{y}(x_k)] \tau_i^* \frac{(x_k - \bar{x}_i)^2}{\sigma_i^3}.$$

If γ is the chosen step length, $e = y_k - \hat{y}(x_k)$ is the actual error at the data point (x_k, y_k) presented, and on-line sequential tuning is performed in the direction of the negative gradient, then the simple gradient technique yields

$$\hat{y}_i := \hat{y}_i + \gamma e \tau_i^*,$$

$$\bar{x}_i := \bar{x}_i + \gamma e [\hat{y}_i - \hat{y}(x_k)] \tau_i^* (x_k - \bar{x}_i) / \sigma_i^2,$$

$$\sigma_i := \sigma_i + \gamma e [\hat{y}_i - \hat{y}(x_k)] \tau_i^* (x_k - \bar{x}_i)^2 / \sigma_i^3.$$

In Fuzzy Logic Toolbox of MATLAB the function `genfis2` can be used for generating a Sugeno fuzzy system by using subtractive clustering, also for the multi-input case.

III. INPUT-OUTPUT SPECIFICATIONS

In order to execute clustering algorithms effectively, it is necessary to extract the main features and relevant data fields from the database. It is useful to derive new variables characterizing fraudulent behavior and build up an abstract space where the function approximation of fraudulence can be accomplished. Based on subjective hypothesis based on the experience of employees on the charge of customer control at the service provider, we have defined 11 input variables which may play the informer when a customer

commits fraudulence. The input following variables are used for the subtractive clustering algorithm:

1. $T_N - T_1$, No. of days between the first (T_1) and last (T_N) known data reading of the consumption meter (*integer*).
2. No. of readings of consumption meter (*integer*).
3. Average time interval (in days) between data readings on the consumption meter (*double*)

$$\frac{\sum_{i=1}^{N-1} (T_{i+1} - T_i)}{N - 1}$$

where N is the Identification number of the consumption meter.

4. Average daily consumption between the first and last data reading on the consumption meter (*double*)

$$\frac{m_N - m_1}{T_N - T_1}$$

where m_i is the value consumption counter at the i th data reading.

5. Flag for discounted nightly consumption {0,1}.
6. The rate of data readings carried out by the customer in relative to total data readings [0,1].
7. Derived change on the counter of consumption meter (*double*).

$$\sum_{i=3}^N \left(\frac{\frac{m_{i-1} - m_{i-2}}{T_{i-1} - T_{i-2}} - \frac{m_i - m_{i-1}}{T_i - T_{i-1}}}{\frac{m_{i-1} - m_{i-2}}{T_{i-1} - T_{i-2}}} \right)^2$$

Exceptions (*integer*):

- 2: No consumption, i.e $m_i = m_{i-1} \forall i = 2 \dots$
 - 1: Backward counted $\exists i: m_i < m_{i-1}$
 - 0: No sufficient information is available.
8. No. of customer control carried out by the provider related to the consumption meter (*integer*).
 9. No. of jobs at the place of consumption meter not related to control process
 10. No. of failed jobs/ No. of total jobs at the location of consumption meter [0,1].
 11. The identification number of the type of consumption meter (*integer*).

The first input parameter is not a real input parameter and not used for training. The only reason to consider it input variable is that it identifies the customer when the control list is generated. The output variables of fuzzy inference systems are defined in a straightforward way:

1. Reported for violation of contract {0,1}
2. Reported for missing contract {0,1}
3. Reported for measurement error {0,1}
4. Reported for damage {0,1}
5. Reported for permanently closed address {0,1}

6. Reported for any reason {0,1}

Note that the sixth output variable is not independent from the others, as it can be derived as their disjunction.

When parameters are chosen for the algorithm, the discrete value range of some input and all the output variables should be taken into account. For cluster radius, 0.25 was chosen to each input variable.

IV. EVALUATION

Meters in the Customer Registration Database were divided into two sets. If the last check of the given meter was carried out prior to 1/1/2009 (229 909 records), its data was used as training data. Otherwise, its data was included to the set of validation data. In this group, there are 94 827 records. The fuzzy inference system was set up by subtractive clustering on the base of training data and validated on the validation data. Subtractive clustering algorithm provides an ordered list of the customers called Control List. Control list contains all the 94 827 records from the second data group (validation data). Lower position in the list represents higher probability for the fraudulence at the customer's address. Naturally, electricity provider intends to check the customers on the top of the Control List. Since visiting a customer's address has a cost, electricity provider has to decide how many customers should be checked. Fig. 1.-Fig. 6 show some correlations which might help the service provider to choose the customers for actual check. Each figure consists of two plots. The upper plot shows the rate of hits as a function of the rate of controlled population. For example, if the rate of hits is 10 when the rate of controlled population is 0.5 then it means that checking the first 0.5% of the customers (1150 customers) from the Control List reveals 10% of the checked population as fraudulent customers. It is emphasized that in the first part, we carried out a simulation to test the algorithm. It means that it was actually known how many fraudulent customers are in the validation data. On the lower part of each figure, the rate of revealed cases relative to the total number of fraudulent population is shown with respect to the rate of controlled population (i.e. validation data).

Horizontal line on the upper part of figures shows the rate of reported cases on the entire population. It is seen that the rate of hits (i.e. the discovery rate) fluctuates if the rate of controlled population is small. The reason behind it is that every hit in small group changes significantly the discovery hit. Also, it may occur that discovery rate is zero at the very beginning (until there is no hit). Naturally, discovery rate converges to the horizontal line (rate of reported cases in total) as the rate of controlled population reaches 100%. It can be generally observed around the 0.5% rate of controlled population the rate of hits jumps to a significantly high value, much higher than the horizontal line. From this value, there may be smaller or higher fluctuation but discovery rate converges to the rate of total reported cases (horizontal lines) exclusively from higher value. This phenomenon proves the efficiency of the algorithm since it shows that the performance of the algorithm is always better than a random

Control List. The performance is especially good in the case of Contract violations (Fig. 1.) where the discovery rate is 14% for 100 checked customers, 3,8% for the 500 checked customers. For comparison the rate of total reported case is very small on the validation data (0.114%).

Also, it can be seen from the lower part of Fig. 1. that 50% of the fraudulent cases can be revealed by checking 5% of the controlled population and 80% of the fraudulent cases can be revealed by checking 20% of the controlled population. Similar trends (though not so striking numerical results) are valid for the other categories of reported cases, as well. It is worth remarking that usual discovery rate is around 3-4% (for 100-1000 checks) which relies on the experience of provider's control unit and can be considered an acceptably good result for previous (human) strategies.

Considering the results on Fig. 1.-Fig. 6 and the fact that the rate of fraudulent cases on the validation data was much less (only 0.114%) than 3-4%, the performance of our procedure seems to be quite promising. Based on the simulation result and the cost of visiting customers' address, electricity provider may calculate how many customers worth to be checked from the Control List. Encouraged by simulation data, Hungarian Electricity Provider ELMŰ-ÉMÁSZ Company-Group checked 500 addresses to reveal fraudulent customers. The procedure was carried out by both divisions of the company groups (ELMŰ and ÉMÁSZ). The results are shown by Table I.

Discovery rates were similar to the simulation data which reinforced the efficiency of the algorithm.

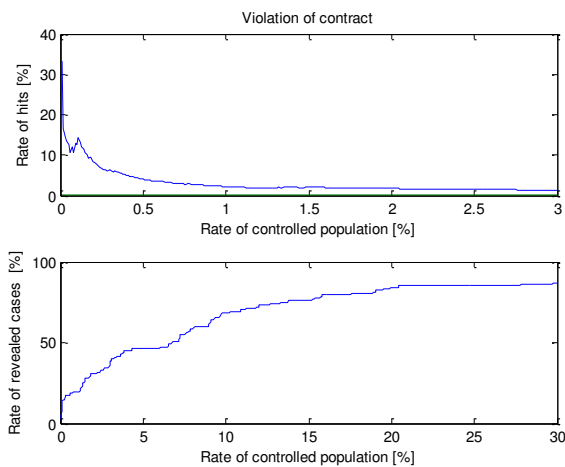


Fig. 1. Rate of hits and rate of revealed cases for violation of contracts.

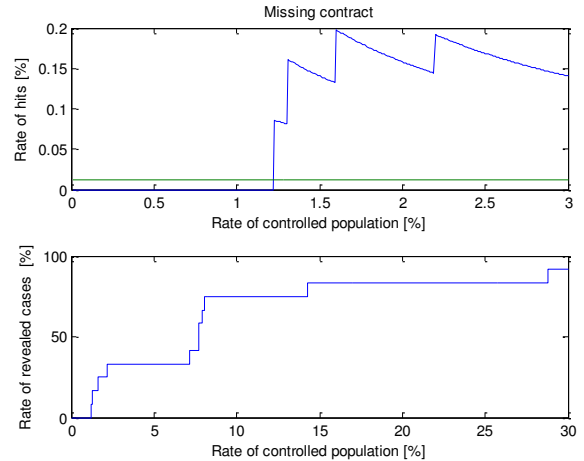


Fig. 2. Rate of hits and rate of revealed cases for missing contracts.

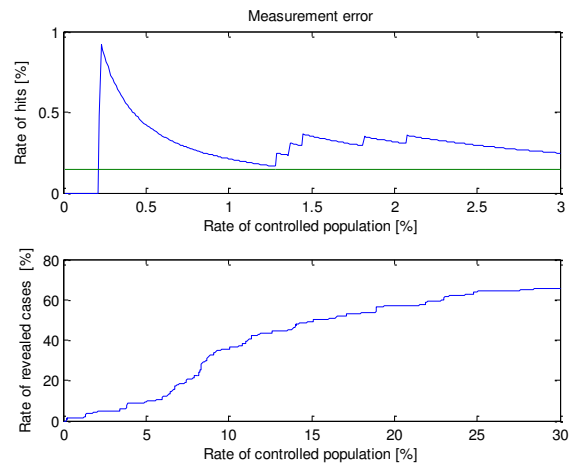


Fig. 3. Rate of hits and rate of revealed cases for measurement errors.

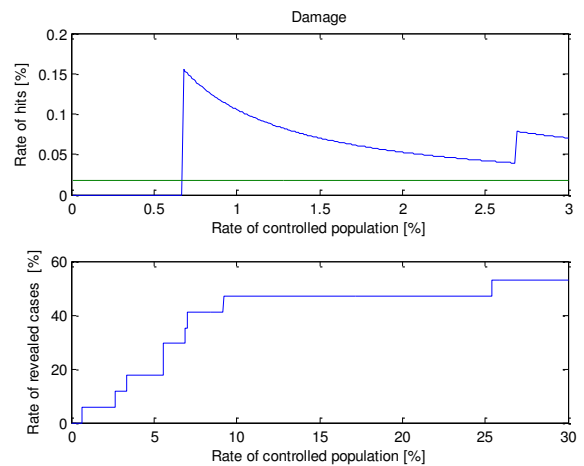


Fig. 4. Rate of hits and rate of revealed cases for damages.

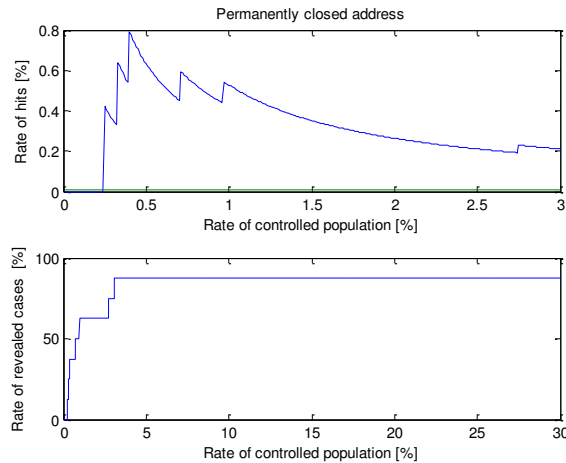


Fig. 5. Rate of hits and rate of revealed cases for permanently closed addresses.

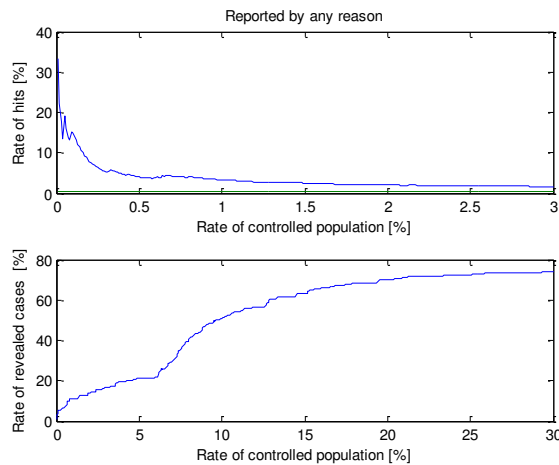


Fig. 6. Rate of hits and rate of revealed cases for any reports

TABLE I

First 50 addresses on the control list are checked.					
	Issued addresses	Visited addresses	Completed checks	Reports	Discovery rate (%)
ELMŰ	19	9	4	1	25
ÉMÁSZ	31	27	26	3	11,5
Total	50	36	30	4	13,3
All the 500 addresses on the control list are checked.					
	Issued addresses	Visited addresses	Completed checks	Reports	Discovery rate (%)
ELMŰ	300	197	99	8	8,1
ÉMÁSZ	200	172	165	7	4,2
Total	500	369	264	15	5,7

V. CONCLUSION

A control list generation method for fraudulence detection has been proposed in this paper. Fraudulence discovery was carried out by a fuzzy system where the rules were determined by subtractive clustering method on the base of training data. After database transformation and input-output

specifications, the algorithm was able to achieve a good performance on validation data set. Hungarian Electricity Provider ELMŰ-ÉMÁSZ Company Group actually checked the offered customer addresses for fraudulence. Realized discovery rate (5-25%) has forcefully exceeded the expectation of 3-4%.

ACKNOWLEDGMENT

Authors thank Tibor Kovács, Zoltán Antal, Balázs Kollár, Anita Borbély, Róbert Bukovics, István Kiss, Botond Kiss, István Kormos, János Rotter, the employees of ELMŰ-ÉMÁSZ Company-Group for their help provided about the Company Group's Customer Registration System and their experiences related to the set of fraudulent customers.

REFERENCES

- [1] I. van Langevelde, A. Philippen, J. Treur, "A Compositional Architecture for Simple Design Formally Specified in DESIRE," *In J. Treur and Th. Wetter, eds., Formal Specification of Complex Reasoning Systems*, Ellis Horwood, New York, 1993.
- [2] A. Th. Schreiber, B. Wielinga, J. Breuker: KADS, "A Principled Approach to Knowledge-Based System Development," Knowledge-Based Systems Vol. 11, Academic Press, London, 1993
- [3] M. Taboada, M. Argüello, J. Des, J. Mira, "Building Knowledge-Based Intelligent Systems by Reusing," *In Innovations in KE, International Series on Advanced Intelligence, Vol. 82, pp. 1-30*, 2003.
- [4] Kovács G., Piétrac L, Kiss B., Niel E., "On the Formalization of Integrating Watchdogs into Supervisory Controller Structures." *In: European Control Conference. Kos, Greece, IEEE, pp. 5522-5529*. 2007
- [5] Lantos, B., "Fuzzy systems and genetic algorithms," Műegyetemi kiadó 2002
- [6] N. Grira, M. Crucianu, N. Boujemaa, "Unsupervised and Semi-supervised Clustering: a Brief Survey," A Review of Machine Learning Techniques for Processing Multimedia Content, Report of the MUSCLE European Network of Excellence (6th Framework Programme), 2005
URL: www.rocq.inria.fr/~crucianu/src/BriefSurveyClustering.pdf
- [7] S. Chiu, "Fuzzy Model Identification Based on Cluster Estimation," *J. of Intelligent & Fuzzy Systems*, Vol. 2, No. 3, 1994.
- [8] Jang, J.-S. R., Sun, C.-T., Mizutani, E., "NeuroFuzzy and Soft Computing – A Computational Approach to Learning and Machine Intelligence," Prentice Hall, 1997
- [9] Azuaje, F., Dubitzky, W., Black, N., Adamson, K., "Discovering Relevance Knowledge in Data: A Growing Cell Structures Approach," *IEEE Transactions on Systems, Man, and Cybernetics Part B: Cybernetics*, Vol. 30, No. 3, June 2000 (pp. 448).
- [10] Lin, C., Lee, C., "Neural Fuzzy Systems," Prentice Hall, NJ, 1996
- [11] Nauck, D., Kruse, R., Klawonn, F., "Foundations of Neuro-Fuzzy Systems," John Wiley & Sons Ltd., NY, 1997. R. Faulhaber, "Design of service systems with priority reservation," *in Conf. Rec. 1995 IEEE Int. Conf. Communications*, pp. 3–8.
- [12] Jang, J.-S. R., "ANFIS: Adaptive-Network-based Fuzzy Inference Systems," *IEEE Transactions on Systems, Man, and Cybernetics*, Vol. 23, No. 3, pp. 665-685, May 1993
- [13] Chiu, S., "Fuzzy Model Identification Based on Cluster Estimation," *Journal of Intelligent & Fuzzy Systems*, Vol. 2, No. 3, Spt. 1994