

A hybrid method for extraction of logical rules from data

Włodzisław Duch, Rafał Adamczak, Krzysztof Grąbczewski and Grzegorz Żal
Department of Computer Methods, Nicholas Copernicus University,
Grudziądzka 5, 87-100 Toruń, Poland.

E-mail: duch,raad,kgrabcze,dzezik@phys.uni.torun.pl

Abstract

A hybrid method for extraction of logical rules from data has been developed. The hybrid method is based on a constrained multi-layer perceptron (C-MLP2LN) neural network for selection of relevant features and extraction of preliminary set of logical rules, followed by a search-based optimization method using global minimization technique. Constraints added to the cost function change the MLP network smoothly into a network performing logical operations. The method is applicable for symbolic and continuous features, finding optimal linguistic variables. Results for several medical and other data sets show that such hybrid technique finds very simple and highly accurate rules, frequently giving results that are more accurate than those obtained by any other classifier. Crisp logical rules are found first, followed by fuzzy rules only if the accuracy of the crisp rules is not satisfactory. Comparison with other rule extraction methods shows superiority of the hybrid approach. The method is also applicable in data mining problems.

Keywords

Neural networks, logical rules, optimization, machine learning.

I. Logical rules - introduction.

Why should one use logical rules if other methods of classification – machine learning, pattern recognition or neural networks – may be easier to use and give better results? Adaptive systems M_W , including the most popular neural network models, are useful classifiers that adjust internal parameters W performing vector mappings from the input to the output space $Y^{(p)} = M_W(X^{(p)})$. Although they may achieve high accuracy of classification the knowledge acquired by such systems is represented in a set of numerical parameters and architectures of networks in an incomprehensible way. Logical rules should be preferred over other methods of classification provided that the set of rules is not too complex and classification accuracy is sufficiently high. Surprisingly, in some applications simple rules proved to be more accurate and were able to generalize better than many machine and neural learning algorithms. Perhaps the main reason for such good performance of logical rules is related to the problem of finding an optimal balance between the flexibility of adaptive models and the danger of overfitting the data. Although Bayesian regularization [1] may help in case of some neural and statistical classification models, logical rules give much better control over the complexity of the data representation.

Many statistical, pattern recognition and machine learning [2] methods of finding logical rules have been designed in the past. Neural networks are also used to extract logical rules and select classification features. Unfortunately systematic comparison of neural and machine learning methods is still missing. Many neural rule extraction methods have recently been reviewed and compared experimentally [3], therefore we will not discuss them here. Neural methods focus on analysis of parameters (weights and biases) of trained networks, trying to achieve high fidelity of performance, i.e. similar results of classification by extracted logical rules and by the original networks. Non-standard form of rules, such as M -of- N rules (M out of N antecedents should be true), fuzzy rules, or decision trees [2] are sometimes useful but in this paper we will consider only standard IF ... THEN propositional rules.

In classification problems propositional rules may take several forms. The most general form is: IF $X \in K^{(i)}$ THEN $\text{Class}(X) = C_i$, i.e. if X belongs to the cluster $K^{(i)}$ than its class is $C_i = \text{Class}(K^{(i)})$, the same as for all vectors in this cluster. If clusters overlap non-zero probability of classification $p(C_i|X; M)$ for several classes is obtained. This approach does not restrict the shapes of clusters used in logical rules, but unless the clusters are visualized in some way (a difficult task in highly dimensional feature spaces) it does not give more understanding of the data than any black box classifier. A popular simplification of the most general form of logical rules is to describe clusters using separable “membership” functions. This leads to fuzzy rules, for example in the form:

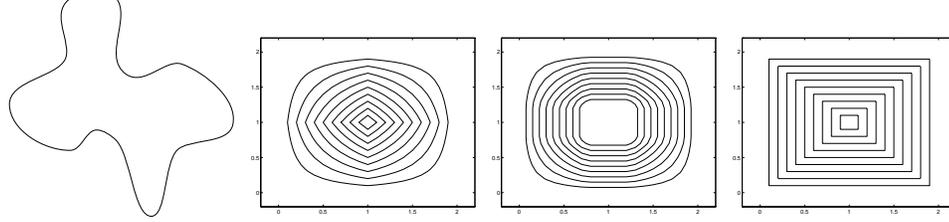


Fig. 1. Shapes of decision borders for general clusters, fuzzy rules (using product of membership function), rough rules (trapezoidal approximation) and logical rules.

$$\mu^{(k)}(X) = \prod_i \mu_i^{(k)}(X_i); \quad p(C_k|X; M) = \frac{\mu^{(k)}(X)}{\sum_i \mu^{(i)}(X)} \quad (1)$$

where $\mu^{(k)}(X)$ is the membership function value for cluster k . Such context-dependent or cluster-dependent membership functions are used in the Feature Space Mapping (FSM) neurofuzzy system [4]. The flexibility of this approach depends on the choice of membership functions. Fuzzy logic classifiers use most frequently a few triangular membership functions per one input feature [5]. These functions do not depend on the region of the input space, providing oval decision borders, similar to Gaussian functions (cf. Fig.1).

Fuzzy rules give decision borders that are not much more flexible than those of crisp rules. What is important is rather the ability to deal with skewed distribution of data by rotating some decision borders. This requires new linguistic variables formed by taking linear combination of input features and the meaning of such rules is sometimes difficult to comprehend.

The rough set theory [6] is used to derive crisp logic propositional rules. For the two-class problems the lower approximation of the data is a set of instances or a region of the feature space containing examples that certainly belong to one class only, while the upper approximation covers all instances which have a chance to belong to this class. In practice the shapes of the boundary between the upper and the lower approximations depends on the indiscernibility (or similarity) relation used. Linear approximation to the boundary region leads to trapezoidal membership functions. The simplest crisp form of logical rules is obtained if trapezoidal membership functions are changed into rectangular functions. Rectangles allow to define logical linguistic variables for each feature by intervals or sets of nominal values.

A fruitful way of looking at logical rules is to treat them as an approximation to the posterior probability of classification $p(C_i|X; M)$, where the model M is composed of the set of rules. Crisp, fuzzy and rough set decision borders are a special case of the FSM neurofuzzy approach [4] based on separable functions used to estimate the classification probability. Although the decision borders of crisp logical rule classification are much simpler than those achievable by neural networks, results are sometimes significantly better, perhaps due to the problem of finding globally optimal solution of the non-linear optimization problem solved with neural classifiers.

We are sure that in all cases, independently of the final classifier used, it is advantageous to extract crisp logical rules. First, in our tests logical rules proved to be highly accurate; second, they are easily understandable by experts in a given domain; third, they may expose the problems with the data itself. This became evident in the case of the ‘‘Hepar dataset’’, collected by dr. H. Wasyluk and co-workers from the Postgraduate Medical Center in Warsaw. The data contained 570 cases described by 119 values of medical tests and other features and were collected over a period of more than 5 years. These cases were divided into 16 classes, corresponding to different types of liver disease (final diagnosis was confirmed by analysis of liver samples under microscope, a procedure that we try to avoid by providing reliable diagnosis system). We have extracted crisp logical rules from this dataset using our MLP2LN method described below, and found that six very simple rules gave 98.5% accuracy. Unfortunately these rules also revealed that the missing attributes in the data were replaced by the averages for a given class, for example, cirrhosis was fully characterized by the rule ‘‘feature₃=4.39’’. Although one may report very good results using crossvalidation tests with this data, such data is of course useless since the averages for a given class are known only after the diagnosis.

In the next section the hybrid rule extraction algorithms used by us are described and ways of estimating accuracy of rules and overcoming their brittleness addressed. Performance of these algorithms is illustrated on a few benchmark and real life problems in the third section. The paper is finished with a short discussion.

II. The hybrid rule extraction methodology

Logical rules require symbolic inputs (linguistic variables). The problem of optimal selection of input features is very important and may be solved in an adaptive way by analysis of the nodes developed by neural networks estimating the probability density, such as the FSM network [4]. Optimal intervals and other adaptive parameters

may be found by maximization of a predictive power of a classifier. Let $P(C_i, C_j|M)$ be the confusion matrix, i.e. the number of instances in which class C_j is predicted when the true class was C_i , given some parameters M . Then for n samples $p(C_i, C_j|M) = P(C_i, C_j|M)/n$ is the probability of (mis)classification. Maximizing the number of correct predictions (called also the “predictive power” of rules):

$$\max_M \text{Tr } P(C_i, C_j|M) \quad (2)$$

over all parameters of the model M , or minimizing the number of wrong predictions (possibly with some cost matrix $R(C_i, C_j)$):

$$\min_M \sum_{i \neq j} R(C_i, C_j) P(C_i, C_j|M) \quad (3)$$

allows to find optimal values of adaptive parameters. Weighted combination of these two terms:

$$E(M) = \lambda \sum_{i \neq j} P(C_i, C_j|M) - \text{Tr } P(C_i, C_j|M) \quad (4)$$

is bounded by $-n$ and should be minimized without constraints. This cost function allows to reduce the number of errors to zero for models M that reject some instances. Since rules discriminate between instances of one class and all other classes one can define a cost function for each rule separately:

$$E_R(M) = \lambda(P_{+-} + P_{-+}) - (P_{++} + P_{--}) \quad (5)$$

and minimize it over parameters M used in the rule \mathcal{R} (+ means here one of the classes, and $-$ means all other classes). The combination $P_{++}/(P_{++} + P_{+-}) \in (0, 1]$ is sometimes called the sensitivity of a rule [7], while $P_{--}/(P_{--} + P_{-+})$ is called the specificity of a rule. Some rule induction methods optimize such combinations of $P_{x,y}$ values. Instead of the number of errors P one may also use the probabilities p .

Estimation of the accuracy of rules is very important, especially in medicine. Tests of classification accuracy should be performed using stratified 10-fold crossvalidation, each time including rule optimization on the training set. Changing the value of λ will produce a series of models with higher and higher confidence of correct classification at the expense of a higher rejection rate. A set of rules may classify some cases at the 100% confidence level; if a new instance is not covered by these rules a set of (usually simpler) rules at lower confidence level is used, until the new instance is classified. In this way a reliable estimation of confidence is possible. The usual procedure is to give only one set of rules, assigning to each rule a confidence factor, for example $c_i = P(C_i, C_i|M) / \sum_j P(C_i, C_j|M)$. This is rather misleading. A rule $\mathcal{R}^{(1)}$ that does not make any errors covers typical instances and its reliability is close to 100%. If a less accurate rule $\mathcal{R}^{(2)}$ is given, for example classifying correctly 90% of instances, the reliability of classification for instances covered by the first rule is still close to 100% and the reliability of classification in the border region (cases covered by $\mathcal{R}^{(2)}$ but not by $\mathcal{R}^{(1)}$) is much less than 90%. Including just these border cases gives much lower confidence factors and since the number of such cases is relatively small the estimate itself has low reliability. A possibility sometimes worth considering is to use a similarity-based classifier (such as the k-NN method or RBF network) to improve accuracy in the border region.

Since multidimensional optimization problems are difficult the number of adaptive parameters of the model M should be reduced first and good starting values provided. Neural-based methods are suitable to select relevant features and provide initial parameters for linguistic variables. Logical (linguistic) input variables s_k for continuous input data components x_i are obtained by dividing the data in distinct (for crisp logic) sets: IF ($x_i \in X_j$) THEN ($s_k = T$). For example, $s_k = s$ may designate the fact that the feature s_k is small, and $s_k = \neg s$ that it is not small. Each quantized feature s will have two or more values, for example represented by the vector $V_{s_1} = (+1, -1, -1\dots)$ for the first value, $V_{s_2} = (-1, +1, -1\dots)$ for the second value etc.

Initial values of the intervals for linguistic variables may sometimes be determined by the analysis of histograms. FSM, our density estimation neurofuzzy network, is initialized using simple clusterization methods [8], for example dendrogram analysis of the input data vectors (for very large datasets rounding of continuous values to lower precision is applied first and then duplicate data vectors are sorted out to reduce the number of vectors). Rectangular functions may be used in FSM for a direct extraction of logical rules. MLP (multilayered perceptron) neural models may find linguistic variables by training a special neural layer (cf. [9]) of L-units. Each L-unit is a combination of sigmoidal functions realizing a function $L(x_i; b_i, b'_i) = \sigma(x_i - b_i) - \sigma(x_i - b'_i)$ parameterized by two biases determining the interval, in which this function has non-zero value. The gain of all sigmoidal functions $\sigma(x)$ is slowly increased during learning. In this process fuzzy membership functions (“soft trapezoids”) are transformed into rectangular functions [4], [9]. Similar smooth transformation is used in the FSM network using biradial transfer functions, which are additionally parameterized combinations of products of $L(x_i; b_i, b'_i)$ functions [10].

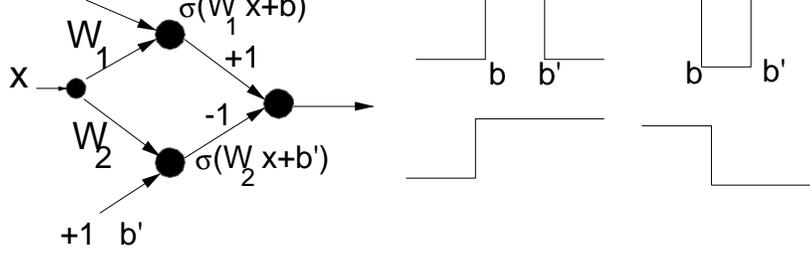


Fig. 2. L-units, or pairs of neurons with constrained weights, used for determination of linguistic variables.

Outputs of L-units $L(x_i; b_i, b'_i)$ are combined and filtered through another sigmoid $\sigma(\sum_{ij} L(x_i; b_{ij}, b'_{ij}))$. Alternatively, a product form $\prod_{ij} L(x_i; b_{ij}, b'_{ij})$ is used. The gain of the sigmoidal functions $\sigma(x)$ is slowly increased during learning and the fuzzy rules are transformed into crisp logical rules, i.e. complex decision regions are transformed into simpler, hypercuboidal decision regions.

To extract rules from trained MLP network one should transform it smoothly into something resembling a logical network (LN). This general idea, called MLP2LN [11], may be realized in several ways. Skeletonization of a large MLP network is the method of choice if our goal is to find logical rules for an already trained network. Otherwise it is simpler to start from a single neuron and construct the logical network using the training data. Since interpretation of the activation of the MLP network nodes is not easy [12] a smooth transition from MLP to a logical-type of network performing similar functions is advocated. This is achieved by: a) gradually increasing the slope of sigmoidal functions to obtain crisp decision regions; b) simplifying the network structure by inducing the weight decay through a penalty term; c) enforcing the integer weight values 0 and ± 1 , interpreted as 0 = irrelevant input, +1 = positive and -1 = negative evidence. These objectives are achieved by adding two additional terms to the standard mean square error function $E_0(W)$:

$$E(W) = E_0(W) + \frac{\lambda_1}{2} \sum_{i,j} W_{ij}^2 + \frac{\lambda_2}{2} \sum_{i,j} W_{ij}^2 (W_{ij} - 1)^2 (W_{ij} + 1)^2 \quad (6)$$

The first term, scaled by the λ_1 hyperparameter, encourages weight decay, leading to skeletonization of the network and elimination of irrelevant features. The second term, scaled by λ_2 , forces the remaining weights to approach ± 1 , facilitating easy logical interpretation of the network function. In the backpropagation training algorithm these new terms lead to the additional change of weights: $\lambda_1 W_{ij} + \lambda_2 W_{ij} (W_{ij}^2 - 1)(3W_{ij}^2 - 1)$.

Instead of the 6-th order term in the cost functions lower order terms may be used:

$$\begin{aligned} & |W_{ij}| |W_{ij}^2 - 1| \text{ cubic} \\ & |W_{ij}| + |W_{ij}^2 - 1| \text{ quadratic} \\ & \sum_{k=-1}^{+1} |W_{ij} + k| - |W_{ij} - \frac{1}{2}| - |W_{ij} + \frac{1}{2}| - 1 \quad \text{linear} \end{aligned} \quad (7)$$

Introduction of integer weights may also be justified from the Bayesian perspective [1]. The cost function specifies our prior knowledge about the probability distribution $P(W|M)$ of the weights in our model M . For classification task when crisp logical decisions are required the prior probability of the weight values should include not only small weights but also large positive and negative weights distributed around ± 1 , for example:

$$P(W|M) = Z(\alpha)^{-1} e^{-\alpha E(W|M)} \propto \prod_{ij} e^{-\alpha_1 W_{ij}^2} \prod_{ij} e^{-\alpha_2 |W_{ij}^2 - 1|} \quad (8)$$

where the parameters α play similar role for probabilities as the parameters λ for the cost function. Using alternative cost functions amounts to different priors for regularization, for example Laplace instead of Gaussian prior. Initial knowledge about the problem may also be inserted directly into the network structure, defining initial conditions modified further in view of the incoming data. Since the final network structure becomes quite simple insertion of partially correct rules to be refined by the learning process is quite straightforward.

Although constraints Eq. (6) do not change the MLP exactly into a logical network they are sufficient to facilitate logical interpretation of the final network function. MLPs are trained using relatively larger λ_1 value at the beginning, followed by large λ_2 value near the end of the training process. The slopes of sigmoidal functions are gradually increased to obtain sharp decision boundaries. Rules \mathcal{R}_k implemented by trained networks are obtained in the form of logical conditions by considering contributions of inputs for each linguistic variable s , represented by a vector V_s . Contribution of variable s to the activation is equal to the dot product $V_s \cdot W_s$ of the subset W_s of the weight vector corresponding to V_s . A combination of linguistic variables activating the hidden neuron above the threshold is a logical rule of the form: $\mathcal{R} = (s_1 \wedge \neg s_2 \wedge \dots \wedge s_k)$. It should be stressed that the only way to change MLP into a logical network is by increasing the slope of sigmoid functions to infinity, changing them into

a step-functions. Such a process is difficult since a very steep sigmoid functions leads to the non-zero gradients only in small regions of the feature space (and thus the number of vectors contributing to the learning process goes to zero). Therefore it is necessary to stop learning at some point and optimize the intervals of the extracted rules calculating errors using the rules instead of MLP network.

In the **constructive version** of the MLP2LN approach usually one hidden neuron per output class is created and training proceeds until the modified cost function reaches minimum. The weights and the threshold obtained are then analyzed and the first group of logical rules is found, covering the most common input-output relations. The input data that are correctly handled by the first group of neurons will not contribute to the error function, therefore the weights of these neurons are kept frozen during further training. This is equivalent to training one neuron (per class) at a time on the remaining data, although sometimes training two or more neurons may lead to faster convergence. This procedure is repeated until all the data are correctly classified, weights analyzed and a set of rules $\mathcal{R}_1 \vee \mathcal{R}_2 \dots \vee \mathcal{R}_n$ is found for each output class or until the number of rules starts to grow rapidly. The output neuron for a given class is connected to the hidden neurons created for that class – in simple cases only one neuron may be sufficient to learn all instances, becoming an output neuron rather than a hidden neuron. Output neurons performing summation of the incoming signals are linear.

Since each time only one neuron per class is trained the C-MLP2LN training is fast. The network repeatedly grows when new neurons are added, and then shrinks when connections are deleted. Since the first neuron for a given class is trained on all data for that class the rules it learns are most general, covering largest number of instances. Therefore rules obtained by this algorithm are ordered, starting with rules that are used most often and ending with rules that handle only a few cases. The final solution may be presented as a set of rules or as a network of nodes performing logical functions. However, some rules obtained from analysis of the network may involve spurious conditions and therefore in the final step should be carefully analyzed and simplified. Once the rules are derived linguistic variables are optimized by minimization of the cost function Eq. (4).

Logical rules, similarly as any other classification systems, may become brittle if the decision borders are placed too close to the data vectors instead of being placed between the clusters. **The brittleness problem** is solved either at the optimization stage by selecting the middle values of the intervals for which best performance is obtained or, in a more general way, by adding noise to the data [13]. So far we have used the simplest method: starting from the initial values of the intervals obtained from MLP network the range of each parameter minimizing the error function is successively determined, the middle value selected, and after the last parameter has been optimized the process is repeated iteratively until the minimum of the error function is reached. Although this process does not guarantee an absolute minimum of the error function it goes beyond simple gradient optimization since the one-dimensional interval optimizations are not restricted to local minima, but are performed over the whole range of admissible parameter values. If several disjointed intervals give the same error a search tree is generated and all branches explored. The center of the cuboid in the parameter space is taken as the final estimation of the adaptive parameters of logical rules.

III. Illustrative applications.

Datasets for benchmark applications were taken from the UCI machine learning repository [14]. Application of the constructive MLP2LN approach to the classical Iris dataset was already presented in detail [15], therefore only new aspects related to the hybrid method are discussed here. The Iris data has 150 vectors evenly distributed in three iris-setosa, iris-versicolor and iris-virginica classes. Each vector has four features: sepal length x_1 and width x_2 , and petal length x_3 and width x_4 (all in cm). Analysis of smoothed histograms (assuming Gaussian width for each value) of the individual features for each class provides initial linguistic variables. Assuming at most 3 linguistic variables per input feature a network with 12 binary inputs equal to ± 1 (features present or absent) is constructed. For example, the medium value of a single feature is coded by $(-1, +1, -1)$ vector. For the Iris dataset a single neuron per one class was sufficient to train the network, therefore the final network structure (Fig. 3) has 12 input nodes and 3 output nodes (hidden nodes are only needed when more than one neuron is necessary to cover all the rules for a given class).

The constraint hyperparameters were increased from $\lambda = 0.001$ at the beginning of the training to about $\lambda = 0.1$ near the end, with larger λ_1 to enforce weight decay at the beginning and larger λ_2 to enforce integer weights near the end of training. On average the network needed about 1000 epochs for convergence. The final weights are taken to be exactly ± 1 or 0 while the final value of the slopes of sigmoids reaches 300. Using L-units with initial linguistic variables provided by the histograms very simple networks were created, with non-zero weights for only one attribute, petal length x_3 . Two rules, giving overall 95.3% accuracy (7 errors) are obtained:

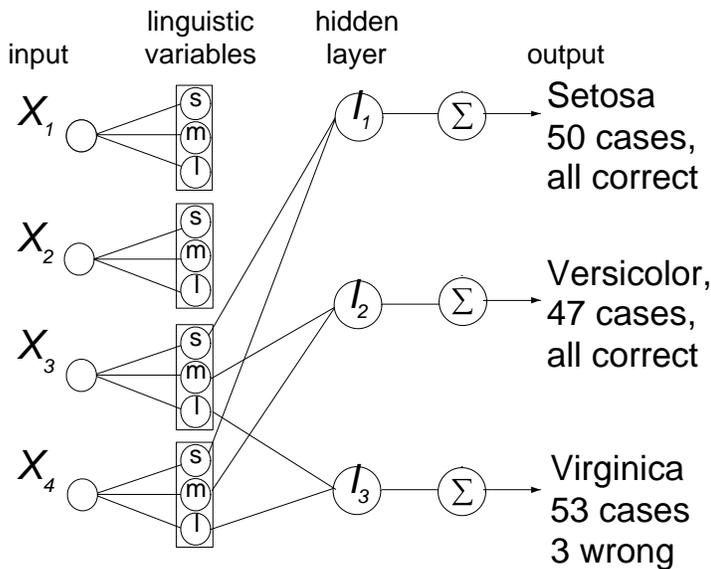


Fig. 3. Final structure of the network for the Iris problem.

$\mathcal{R}_1^{(1)}$: IF $x_3 < 2.5$ THEN Iris-Setosa (100%)
 $\mathcal{R}_2^{(1)}$: IF $x_3 > 4.8$ THEN Iris-Virginica (92%)
 $\mathcal{R}_3^{(1)}$: IF ELSE THEN Iris-Versicolor (94%)

Lowering the final hyperparameters leads to the following weights and thresholds (only the signs of the weights are written):

Setosa	(0,0,0	0,0,0	+,0,0	+,0,0)	$\theta = 1$
Versicolor	(0,0,0	0,0,0	0,+,-	0,+,-)	$\theta = 3$
Virginica	(0,0,0	0,0,0	-,+,-	-,+,-)	$\theta = 2$

Interpretation of these weights and the resulting network function (Fig. 3) is quite simple. Only two features, x_3 and x_4 , are relevant and a single rule per class is found:

IF $(x_3 < 2.9 \vee x_4 < 0.9)$ THEN Iris-Setosa (100%)
 IF $(x_3 \in [2.9, 4.95] \wedge x_4 \in [0.9, 1.65])$ THEN Iris-Versicolor (100%)
 IF $(x_3 > 4.95) \vee (x_4 > 1.65)$ THEN Iris-Virginica (94%)

This $\mathcal{R}^{(2)}$ level set of rules allows for correct classification of 147 vectors, achieving overall 98.0% accuracy. However, the first two rules have 100% reliability while all errors are due to the third rule, covering 53 cases. Decreasing constraint hyperparameters further allows to replace one of these rules by four rules, with a total of three attributes and 11 antecedents, necessary to classify correctly a single additional vector, a clear indication that overfitting occurs. Rules at the $\mathcal{R}^{(3)}$ level that are 100% reliable reject 11 vectors, 8 virginica and 3 versicolor:

IF $x_3 < 2.9$ THEN Iris-Setosa (100%)
 IF $x_3 \in [2.9, 4.9] \wedge x_4 < 1.7$ THEN Iris-Versicolor (100%)
 IF $x_3 \geq 5.3 \vee x_4 \geq 1.9$ THEN Iris-Virginica (100%)

The reliability of classification for the 11 vectors in the border region ($\mathcal{R}^{(2)}$ minus $\mathcal{R}^{(3)}$) is rather low: with $p=8/11$ they should be assigned to the virginica class and with $p=3/11$ to the versicolor class. It is possible to generate more specific rules, including more features, just for the border region. Unfortunately since there are only a few vectors there reliability of these rules will be low and for the Iris data set a simple rule that would correctly separate the 3 versicolor vectors in the border region does not exist. It is also possible to use in this region similarity-based classification system, such as the k-nearest neighbors method, but for this small dataset we do not expect any real improvement since the true probability distributions of leave's sizes for the two classes of iris flowers clearly overlap.

In the **mushroom problem** [3], [14] the database consists of 8124 vectors, each with 22 discrete attributes with up to 10 different values. 51.8% of the cases represent edible, and the rest non-edible (mostly poisonous) mushrooms.

A single neuron is capable of learning all the training samples (the problem is linearly separable), but the resulting network has many nonzero weights and is difficult to analyze from the logical point of view. Using the constructive MLP2LN algorithm with the cost function Eq. (6) the following disjunctive rules for poisonous mushrooms have been discovered:

\mathcal{R}_1	odor= \neg (almond \vee anise \vee none)	(100%)
\mathcal{R}_2	spore-print-color=green	(100%)
\mathcal{R}_3	odor=none \wedge stalk-surface-below-ring=scaly \wedge (stalk-color-above-ring= \neg brown)	(100%)
\mathcal{R}_4	habitat=leaves \wedge cap-color=white	(100%)

Rule \mathcal{R}_1 misses 120 poisonous cases (giving overall 98.52% accuracy), adding rule \mathcal{R}_2 leaves 48 errors (99.41% accuracy), adding third rule leaves only 8 errors (99.90% accuracy), and rules $\mathcal{R}_1 - \mathcal{R}_4$ classify all poisonous cases correctly. The first two rules are realized by one neuron, adding a second neuron and training it on the remaining cases generates two additional rules, \mathcal{R}_3 handling 40 cases and \mathcal{R}_4 handling only 8 cases. In \mathcal{R}_4 "habitat=leaves" may also be replaced by "population=clustered". This is the simplest systematic logical description (some of these rules have probably been also found by the RULEX and TREX algorithms [3]) of the mushroom dataset that we know of and therefore should be used as benchmark for other rule extraction methods.

We have also solved the three monk problems [16]. For the Monks 1 problem one additional neuron handling exceptions has been generated, giving a total of 4 rules and one exception and classifying the data without any errors. In the Monk 2 problem perfect classification with 16 rules and 8 exceptions extracted (exceptions are additional rules handling patterns that are not recognized properly by the rules) from the resulting network has been achieved. The number of atomic formulae which compose them is 132. In the Monk 3 problem, although the training data for this problem is corrupted by 5% noise it is still possible to obtain 100% accuracy [3]. The whole logical system for this case contains 33 atomic formulae.

IV. Applications to medical data

To facilitate comparison with results obtained by several classification methods we have selected three well-known medical datasets obtained from the UCI repository [14].

A. Wisconsin breast cancer data.

The Wisconsin cancer dataset [17] contains 699 instances, with 458 benign (65.5%) and 241 (34.5%) malignant cases. Each instance is described by the case number, 9 attributes with integer value in the range 1-10 (for example, feature f_2 is "clump thickness" and f_8 is "bland chromatin") and a binary class label. For 16 instances one attribute is missing. This data has been analyzed in a number of papers.

Initially 5 rules from the C-MLP2LN procedure were obtained for malignant cases, with benign cases covered by the ELSE condition, with overall accuracy of 96%.

$\mathcal{R}1$:	$f_2 < 6 \wedge f_4 < 4 \wedge f_7 < 2 \wedge f_8 < 5$	(100%)
$\mathcal{R}2$:	$f_2 < 6 \wedge f_5 < 4 \wedge f_7 < 2 \wedge f_8 < 5$	(100%)
$\mathcal{R}3$:	$f_2 < 6 \wedge f_4 < 4 \wedge f_5 < 4 \wedge f_7 < 2$	(100%)
$\mathcal{R}4$:	$f_2 \in [6, 8] \wedge f_4 < 4 \wedge f_5 < 4 \wedge f_7 < 2 \wedge f_8 < 5$	(100%)
$\mathcal{R}5$:	$f_2 < 6 \wedge f_4 < 4 \wedge f_5 < 4 \wedge f_7 \in [2, 7] \wedge f_8 < 5$	(92.3%)

The first 4 rules achieve 100% accuracy (on the whole dataset), the last rule covers only 39 cases, including 3 errors. The confusion matrix is: $P = \begin{pmatrix} 238 & 3 \\ 25 & 433 \end{pmatrix}$, i.e. there are 3 benign cases wrongly classified as malignant and 25 malignant cases wrongly classified as benign, giving overall accuracy of 96%. The simplest form of rules obtained from optimization includes only two rules for malignant class:

$$\text{IF } f_2 \geq 7 \vee f_7 \geq 6 \text{ THEN malignant} \quad (95.6\%)$$

These rules cover 215 malignant cases and 10 benign cases, achieving overall accuracy (including ELSE condition) of 94.9%. More accurate set of rules has been obtained from optimization procedure. Only 1 benign vector is classified as malignant ($\mathcal{R}1$ and $\mathcal{R}5$, the same vector). The ELSE condition for the benign class makes 6 errors, giving 99.00% overall accuracy of this set of rules:

TABLE I

Results from the 10-fold crossvalidation for the Wisconsin breast cancer dataset.

Method	Accuracy %	Reference
IncNet	97.1	[18]
3-NN, Manhattan	97.0±0.12	(our own results)
Fisher LDA	96.8	[19]
MLP+backprop	96.7	[19]
LVQ	96.6	[19]
Bayes (pairwise dependent)	96.6	[19]
FSM (our results)	96.5	(our own results)
Naive Bayes	96.4	[19]
Linear Discriminant Analysis	96.0	[19]
CART (decision tree)	94.2	[19]
LFC, ASI, ASR decision trees	94.4-95.6	[19]
Quadratic Discriminant Analysis	34.5	[19]

$$\begin{aligned}
\mathcal{R}_1) \quad & f_2 < 6 \wedge f_4 < 3 \wedge f_8 < 8 && (99.8\%) \\
\mathcal{R}_2) \quad & f_2 < 9 \wedge f_5 < 4 \wedge f_7 < 2 \wedge f_8 < 5 && (100\%) \\
\mathcal{R}_3) \quad & f_2 < 10 \wedge f_4 < 4 \wedge f_5 < 4 \wedge f_7 < 3 && (100\%) \\
\mathcal{R}_4) \quad & f_2 < 7 \wedge f_4 < 9 \wedge f_5 < 3 \wedge f_7 \in [4, 9] \wedge f_8 < 4 && (100\%) \\
\mathcal{R}_5) \quad & f_2 \in [3, 4] \wedge f_4 < 9 \wedge f_5 < 10 \wedge f_7 < 6 \wedge f_8 < 8 && (99.8\%)
\end{aligned}$$

In all cases features f_3 and f_6 (both related to the cell size) were not important and f_2 with f_7 being the most important. Minimizing Eq. (4) we have initially obtained 7 rules (3 for malignant and 4 for benign class) working with 100% reliability but rejecting 79 cases (11.3%). Further optimization preserving 100% reliability leads to a set of rules that reject 51 cases (7.3%) of all vectors. For malignant class these rules are:

$$\begin{aligned}
\mathcal{R}_1) \quad & f_2 < 9 \wedge f_4 < 4 \wedge f_7 < 3 \wedge f_8 < 6 \\
\mathcal{R}_2) \quad & f_2 < 5 \wedge f_5 < 8 \wedge f_7 < 5 \wedge f_8 < 10 \\
\mathcal{R}_3) \quad & f_2 < 4 \wedge f_4 < 2 \wedge f_5 < 3 \wedge f_7 < 7 \\
\mathcal{R}_4) \quad & f_2 < 10 \wedge f_5 < 10 \wedge f_7 \in [1, 5] \wedge f_8 < 2
\end{aligned}$$

For the benign cases initial rules are obtained by negation of the above rules; after optimization the rule becomes: $\neg(\mathcal{R}_5 \vee \mathcal{R}_6 \vee \mathcal{R}_7 \vee \mathcal{R}_8)$, where:

$$\begin{aligned}
\mathcal{R}_5) \quad & f_2 < 8 \wedge f_4 < 5 \wedge f_8 < 4 \\
\mathcal{R}_6) \quad & f_2 < 9 \wedge f_5 < 6 \wedge f_7 < 9 \wedge f_8 < 5 \\
\mathcal{R}_7) \quad & f_2 < 9 \wedge f_4 < 6 \wedge f_5 < 8 \wedge f_7 < 9 \\
\mathcal{R}_8) \quad & f_2 = 6 \wedge f_4 < 10 \wedge f_5 < 10 \wedge f_7 < 2 \wedge f_8 < 9
\end{aligned}$$

B. The Cleveland heart disease data.

The Cleveland heart disease dataset [14] (collected at V.A. Medical Center, Long Beach and Cleveland Clinic Foundation by R. Detrano) contains 303 instances, with 164 healthy (54.1%) instances, the rest are heart disease instances of various severity. While the database has 76 raw attributes, only 13 of them are actually used in machine learning tests, including 6 continuous features and 4 nominal values. There are many missing values of the attributes. Results obtained with various methods for this data set are collected in Table 2.

After some simplifications the derived rules obtained by the C-MLP2LN approach are:

$$\begin{aligned}
\mathcal{R}1: \quad & \text{thal}=0 \vee \text{thal}=1 \wedge \text{ca}=0.0 && (88.5\%) \\
\mathcal{R}2: \quad & \text{thal}=0 \wedge (\text{ca}=1.0 \vee \text{ca}=3.0) \wedge \text{cp} \neq 2 && (81.0\%) \\
\mathcal{R}3: \quad & (\text{thal}=0 \vee \text{ca}=0.0) \wedge \text{cp} \neq 2 && (85.2\%)
\end{aligned}$$

These rules give 85.5% correct answers on the whole set and compare favorable with the accuracy of other classifiers (Table 2). After further optimization two rules for the first class (healthy) are obtained:

$$\begin{aligned}
\mathcal{R}_1) \quad & (\text{thal}=0 \vee \text{thal}=1) \wedge \text{ca}=0.0 && (88.5\%) \\
\mathcal{R}_2) \quad & (\text{thal}=0 \vee \text{ca}=0.0) \wedge \text{cp} \neq 2 && (85.2\%)
\end{aligned}$$

The ELSE for the second class has 89.2% reliability.

These rules give 85.5% correct answers on the whole set and compare favorable with the accuracy of other classifiers (Table 2).

TABLE II
Results from the 10-fold crossvalidation for the Cleveland heart disease dataset.

Method	Accuracy %	Reference
IncNet	90.0	[18]
Linear Discriminant Analysis	84.5	[19]
Fisher LDA	84.2	[19]
FSM - Feature Space Mapping	84.0	our own results
Naive Bayes	83.4	[19]
Bayes (pairwise dependent)	83.1	[19]
LVQ	82.9	[19]
k-NN, k=27, Manhattan	82.8±0.6	our own results
MLP+backprop	81.3	[19]
CART (decision tree)	80.8	[19]
Quadratic Discriminant Analysis	75.4	[19]
LFC, ASI, ASR decision trees	74.4-78.4	[19]

C. The hypothyroid data.

This is a somewhat larger dataset [14], with 3772 cases for training, 3428 cases for testing, 22 attributes (15 binary, 6 continuous), and 3 classes: primary hypothyroid, compensated hypothyroid and normal (no hypothyroid). The class distribution in the training set is 93, 191, 3488 vectors and in the test set 73, 177, 3178. Initially 4 rules were found, with 99.68% accuracy on the training set and 99.07% error on the test set. For the first class two rules are sufficient (all values of continuous features are multiplied here by 1000):

$$\mathcal{R}_1: \text{TSH} \geq 29 \wedge \text{FTI} < 63$$

$$\mathcal{R}_2: \text{TSH} \in [6.1, 29] \wedge \text{FTI} < 63 \wedge \text{T3} < 20$$

For the second class one rule is created:

$$\mathcal{R}_3: \text{FTI} \in [63, 180] \wedge \text{TSH} \geq 6.1 \wedge \text{on thyroxine=no} \wedge \text{surgery=no}$$

The third classe is covered by the ELSE rule. After further optimization these rules become:

$\mathcal{R}_1:$	$\text{TSH} \geq 30.5 \wedge \text{FTI} < 64.2$	primary hypothyroid	(97.06%)
$\mathcal{R}_2:$	$\text{TSH} \in [6.2, 29.5] \wedge \text{FTI} < 64.1 \wedge \text{T3} < 23.2$	primary hypothyroid	(100%)
$\mathcal{R}_3:$	$\text{TSH} \geq 6.0 \wedge \text{FTI} \in [64.3, 186.7] \wedge \text{TT4} < 148.5 \wedge$ $\text{on thyroxine=no} \wedge \text{surgery=no}$	compensated hypothyroid	(98.96%)
$\mathcal{R}_4:$	ELSE	THEN no hypothyroid	(100%)

These rules make only 4 errors on the training set (99.89%) and 22 errors on the test set (99.36%). They are similar to those found using heuristic version of PVM method by Weiss and Kapouleas [20]. The differences among PVM, CART and C-MLP2LN are for this dataset rather small (Table 3), but other methods, such as well-optimized MLP (including genetic optimization of network architecture) or cascade correlation classifiers, give results that are significantly worse. Poor results of k-NN are especially worth noting, showing that in this case, despite large amount of reference vectors, similarity-based methods are not competitive.

V. Discussion

A new methodology for extraction of logical rules from data has been presented. Neural networks - either density estimation (FSM) or constrained multilayered perceptrons (MLPs) - are used to obtain initial sets of rules. FSM is trained either directly with the rectangular functions or making a smooth transition from biradial functions [10] or trapezoidal functions to rectangular functions. MLPs are trained with constraints that change them into networks processing logical functions, either by simplification of typical MLPs or by incremental construction of networks performing logical functions. In this paper we have used mostly the C-MLP2LN constructive method since it requires less experimentation with various network structures.

TABLE III
Results for the hypothyroid dataset.

Method	% train	% test
k-NN [20]	–	95.27
Bayes [20]	97.03	96.06
MLP+genetic optimization [21]	99.60	98.45
Cascade correl. [21]	100.00	98.48
PVM [20]	99.79	99.33
CART [20]	99.79	99.36
C-MLP2LN	99.89	99.36

The method of successive regularizations introduced by Ishikawa [22] has several features in common with our MLP2LN approach and is capable (although requiring more effort to select the initial network architecture) of producing similar results. Specific form of the cost function as well as the C-MLP2LN constructive algorithm in which neurons are added and then connections are deleted seems to be rather different from other algorithms used for logical rule extraction so far [3]. However, in our experience neural and machine learning methods may serve only for feature selection and initialization of rules performed at the optimization stage based on global minimization (or search) procedures. Optimization leads to rules that are more accurate and simpler, providing in addition sets of rules with different reliability.

Using this hybrid method the simplest logical description for several benchmark problems (Iris, mushroom) has been found. Very good solutions were obtained for the three monk problems. For many medical datasets (only 3 were shown here) very simple and highly accurate results were obtained. It is not quite clear why logical rules work so well, for example in the hypothyroid or the Wisconsin breast cancer case obtaining accuracy which is better than that of any other classifier. One possible explanation for the medical data is that the classes labeled "sick" or "healthy" have really fuzzy character. If the doctors are forced to make yes-no diagnosis they may fit the results of tests to specific intervals, implicitly using crisp logical rules. Another reason may be related to the difficulty of finding globally optimal solutions of the learning problem by neural networks. Logical rules given in this paper may actually be used to initialize MLPs and if batch training algorithms are used the results should not degrade. This issue requires further investigation.

In any case the approach presented here is ready to be used in real world applications and we are at present applying it to complex medical and psychometric data collected in the Psychological Clinic of our University. A computerized version of the Minnesota Multiphasic Personality Inventory (MMPI) test was used, consisting of 550 questions with 5 possible answers each. MMPI evaluates psychological characteristics reflecting social and personal maladjustment, including psychological dysfunction. Hundreds of books and papers were written on the interpretation of this test (cf. review [23]). The raw MMPI data is used to compute 14 coefficients forming a psychogram. At present we have 1465 cases, each classified into one of 34 types (normal, neurotic, alcoholics, schizophrenic, psychopaths, organic problems, malingers etc.) by an expert psychologist. On average 2.5 logical rules per class were derived, involving between 3 and 7 features. A typical rule has the form: IF $f_7 \in [55, 68] \wedge f_{12} \in [81, 93] \wedge f_{14} \in [49, 56]$ THEN paranoia. After optimization these rules will be used in an expert system and evaluated by clinical psychologists.

Acknowledgment: Support by the Polish Committee for Scientific Research, grant 8 T11F 014 14, is gratefully acknowledged.

References

- [1] D.J. MacKay. A practical Bayesian framework for backpropagation networks, *Neural Comput.* 4 (1992) 448-472
- [2] T. Mitchell, *Machine learning*. McGraw Hill 1997
- [3] R. Andrews, J. Diederich, A.B. Tickle. A Survey and Critique of Techniques for Extracting Rules from Trained Artificial Neural Networks, *Knowledge-Based Systems* 8 (1995) 373-389
- [4] W. Duch and G.H.F. Diercksen, Feature Space Mapping as a universal adaptive system, *Computer Physics Communic.* 87: 341-371, 1995; W. Duch, R. Adamczak and N. Jankowski, New developments in the Feature Space Mapping model, 3rd Conf. on Neural Networks, Kule, Poland, Oct. 1997, pp. 65-70
- [5] N. Kasabov, *Foundations of Neural Networks, Fuzzy Systems and Knowledge Engineering*, MIT Press (1996)
- [6] Z. Pawlak, *Rough sets - theoretical aspects of reasoning about data*, Kluwer Academic Publishers 1991
- [7] S.M. Weiss, C.A. Kulikowski, *Computer systems that learn*, Morgan Kaufman, San Mateo, CA 1990
- [8] W. Duch, R. Adamczak and N. Jankowski, Initialization of adaptive parameters in density networks, 3rd Conf. on Neural Networks, Kule, Oct. 1997, pp. 99-104

- [9] W. Duch, R. Adamczak and K. Grąbczewski, Constrained backpropagation for feature selection and extraction of logical rules, in: *Proc. of "Colloquia in AI"*, Łódź, Poland 1996, p. 163–170
- [10] W. Duch and N. Jankowski, New neural transfer functions. *Applied Math. & Comp. Science* 7 (1997) 639-658
- [11] W. Duch, R. Adamczak and K. Grąbczewski, Extraction of logical rules from backpropagation networks. *Neural Processing Letters* 7, 1-9 (1998)
- [12] J.M. Żurada. *Introduction to Artificial Neural Systems*, West Publishing Company, St Paul, 1992.
- [13] C.M. Bishop, Training with noise is equivalent to Tikhonov regularization, *Neural Comput.* 7, 108-116 (1998)
- [14] C.J. Mertz, P.M. Murphy, UCI repository of machine learning databases, <http://www.ics.uci.edu/pub/machine-learning-data-bases>.
- [15] W. Duch, R. Adamczak, K. Grąbczewski, Extraction of logical rules from training data using backpropagation networks, *The 1st Online Workshop on Soft Computing*, 19-30.Aug.1996; <http://www.bioele.nuee.nagoya-u.ac.jp/wsc1/>, pp. 25-30
- [16] W. Duch, R. Adamczak, K. Grąbczewski, Extraction of crisp logical rules using constrained backpropagation networks, *Int. Conf. on Artificial Neural Networks (ICNN'97)*, Houston, 9-12.6.1997, pp. 2384-2389
- [17] K. P. Bennett, O. L. Mangasarian, Robust linear programming discrimination of two linearly inseparable sets, *Optimization Methods and Software* 1, 1992, 23-34.
- [18] N. Jankowski N and V. Kadirkamanathan, Statistical control of RBF-like networks for classification. In: *7th Int. Conf. on Artificial Neural Networks (ICANN'97)*, Lausanne, Switzerland, 1997, pp. 385–390
- [19] B. Ster and A. Dobnikar, Neural networks in medical diagnosis: Comparison with other methods. In: A. Bulsari et al., eds, *Proc. Int. Conf. EANN'96*, pp. 427-430, 1996.
- [20] S.M. Weiss, I. Kapouleas. An empirical comparison of pattern recognition, neural nets and machine learning classification methods. In: J.W. Shavlik and T.G. Dietterich, *Readings in Machine Learning*, Morgan Kauffman Publ, CA 1990
- [21] W. Schiffman, M. Joost, R. Werner, Comparison of optimized backpropagation algorithms, *Proc. of ESANN'93*, Brussels 1993, pp. 97-104
- [22] W. Duch, R. Adamczak, K. Grąbczewski, M. Ishikawa, H. Ueda, Extraction of crisp logical rules using constrained backpropagation networks - comparison of two new approaches, *Proc. of the European Symposium on Artificial Neural Networks (ESANN'97)*, Bruges 16-18.4.1997, pp. 109-114
- [23] J.N. Butcher, S.V. Rouse, Personality: individual differences and clinical assessment. *Annual Review of Psychology* 47 (1996) 87