

Initializing Neural Networks using Decision Trees

Arunava Banerjee
Department of Computer Science
Rutgers University
New Brunswick, NJ 08903
arunava@cs.rutgers.edu

Abstract

Existing approaches to the inductive learning problem include Symbolic and Connectionist algorithms. While the Symbolic approach is generally found to run significantly faster during learning, the Connectionist algorithms are often more accurate at classifying novel examples in the presence of noisy data. This paper presents a technique that determines the topology and initial weights of a neural network using a decision tree, thus combining both approaches. Experimental results on benchmark real-world datasets indicate that this technique outperforms the above mentioned approaches both in efficiency and accuracy.

1 Introduction

The task of inductive learning from examples is to find an approximate definition for an unknown function $f(\mathbf{x})$, given training examples of the form $\langle x_i, f(x_i) \rangle$. Existing algorithms to solve this problem follow either of two basic approaches. Decision tree methods, such as ID3 (Quinlan, 1986) and CART (Breiman, Friedman, Olshen, and Stone, 1984) construct trees whose leaves are labeled with the predicted classification. The connectionist methods, on the other hand, apply neural network learning algorithms such as

the perceptron algorithm (Rosenblatt, 1958), and the error-backpropagation algorithm (Rumelhart, Hinton and Williams, 1986) that learn through incremental changes of weights in a network consisting of elementary units called neurons.

Comparative studies conducted by two different sets of researchers, (Mooney, Shavlik, Towell, and Gove, 1989) and (Weiss, and Kapouleas, 1989) have shown that whereas the decision tree algorithms run significantly faster during training, the connectionist methods almost always perform better at classifying novel examples in the presence of noisy data.

The studies have also demonstrated that when comparisons are confined strictly to the connectionist approach, the back-propagation algorithm outperforms the perceptron algorithm in classifying real world datasets; the reason being that real world datasets are more often than not linearly inseparable.

Back-propagation however, suffers from a host of drawbacks. First, it is extremely slow to train. Second, and more important, variable parameters like the learning rate, the momentum, and the network topology affect the accuracy of the resulting classifier. Setting these parameters to their respective optimal values is a painstaking process that requires additional time and expertise.

When it comes to solving real-world problems, restrictions on time and accuracy become of utmost importance. This is to say that most applications require accurate real-time solutions to their respective problems. The methods mentioned above either come short of attaining the required accuracy for classifying novel examples, or train far too slowly to be functional in such environments.

This paper describes an algorithm which is much faster than back-propagation, and at the same time matches it in accuracy in classifying novel examples. The key idea is to construct a decision tree, convert it into an equivalent network, and then tune the network by training it over the same dataset for a short period of time.

Many researchers have proposed converting decision trees into equivalent neural networks - see for example (Utgoff, 1989), (Sethi, 1990), (Brent, 1991), and (Cios and Liu, 1992). The difference in our approach lies in the fact that whereas all these methods generate networks that classify as accurately as the input decision trees, our method, aided by subsequent tuning, outperforms the decision tree in its accuracy in classifying new examples.

This paper is structured as follows. Section 2 demonstrates how a neural network may be derived from a decision tree. In addition, it highlights certain other advantages of the proposed method. Section 3 clarifies the algorithm by converting an example decision tree into an equivalent neural network. In section 4 experimental results pertaining to three different domains are presented and the new approach is compared to back-propagation and ID3 in their context. Section 5 discusses future research directions based on some known shortcomings of this approach, and section 6 presents certain concluding remarks.

2 The Conversion Technique

As has already been stated, the central idea of the technique is to apply a decision tree (created by any decision tree algorithm, for example ID3) to initialize the neural network. This network is subsequently trained using a connectionist method on the dataset to achieve an improved predictive accuracy. In this paper, C4.5 has been chosen as a proto-typical representative of the decision tree approach. Similarly, Back-propagation (**Bp**) has been chosen as a representative of the connectionist approach. Furthermore, in order to calibrate the method objectively, all experimental comparisons have been made against C4.5 and **Bp**.

To begin with however, we delve into certain important characteristics of the new technique.

It has three major advantages. First, the initial network performs almost as well as the decision tree on which it is based. In other words, even before the process of subsequent training starts, the network is a considerably accurate classifier of the dataset. Second, the procedure that converts the decision tree into a neural network also specifies the network topology. This eliminates the guess work that goes into the creation of a network topology. Third, as will be demonstrated through the experiments, this technique fares better than the best network classifier and decision tree, even when the values for momentum and learning rates are frozen across datasets. Consequently, one could eliminate the process of repeated training required to find the optimal values for such parameters, and at the same time retain the accuracy of the classifier.

An informal sketch of the technique is presented next. This is followed

by the actual algorithm.

2.1 Informal Sketch

The technique presented in this paper converts uni-variate decision trees into equivalent neural networks. The network that is created has three layers (two hidden layers and one output layer). The two hidden layers are called the *literal* layer and the *conjunction* layer. The output layer is also called the *disjunction* layer.

Any data point is assumed to be an element of R^n tagged with a concept class that is a member of the set $\{1, 2, \dots, M\}$. In other words, each example is assumed to be a vector of length $(n + 1)$ that comprises of n real valued attributes and an integral class that takes a value between 1 and M .

The network generated by the technique has exactly n input nodes and M output nodes. The n input nodes correspond to the n attributes in an example, and the ordinality of the node in the output layer that has the maximum activation corresponds to the classification of the example.

The technique first creates a decision tree by applying C4.5 on the dataset. This is followed by a phase of the creation of the actual network. Corresponding to each node in the decision tree, the technique requires the construction of a node in the literal layer that replicates the decision of that node in the tree.

The next phase involves the creation of nodes that correspond to branches in the tree. This is accomplished by assigning one node in the conjunction layer to each branch such that the node replicates the decision of that branch in the tree. The penultimate phase corresponds to the act of grouping together all branches that are assigned to the same output class. This is achieved by allocating one node each in the output (disjunction) layer to a class, and by associating with it all nodes in the conjunction layer that correspond to the branches for that class.

In the final phase, a set of weak edges is superimposed on the network. These edges connect each node in a layer to every node in the previous and the next layer that remain disconnected from it after the previous step.

2.2 The Formal Algorithm

The validity of the technique is based on the assumption that each node in the generated network performs a hyperplane test of the form

$$b + \sum_{j=1}^n w_j x_j > 0? \quad (1)$$

and replies with a zero or one depending on whether the test succeeds. Here, b corresponds to the bias of the node, w_j corresponds to the weight on the input edge j , and x_j to the value of the input on edge j . In order to facilitate the subsequent training using **Bp**, the hard threshold is softened by replacing it with a sigmoid of the form

$$f(x) = \frac{1}{1 + e^{k \cdot (b-x)}} \quad (2)$$

where k characterizes the softness ¹ of the threshold.

The algorithm is as follows.

Begin

- 1 Initialize variables σ and β to 5.0 and 0.025 respectively.
- 2 Run C4.5 on the training dataset to generate a decision tree.
- 3 Traverse the decision tree to create a **dnf** (disjunctive normal form) formula for each class.
- 4 Eliminate all redundant literals from each disjunct.
- 5 For each distinct literal of the form $\langle attrib \rangle > \langle value \rangle$, ² create a hidden unit in the literal layer with a bias of $-\sigma * \langle value \rangle$. Connect it to the input unit corresponding to $\langle attrib \rangle$ with a weight of σ . Connect it to all other input units with weights $+\beta$ or $-\beta$ with equal probabilities.
- 6 For each literal of the form $\langle attrib \rangle < \langle value \rangle$, repeat step 5 with the signs for the bias and weights inverted.

¹The smaller the value of k , the softer is the threshold

²Literals of the form $\langle attrib \rangle = \langle value \rangle$ can be converted to $\langle attrib \rangle > \langle value - \delta \rangle \wedge \langle attrib \rangle < \langle value + \delta \rangle$

- 7 For each disjunct in a class, create a new hidden unit in the conjunction layer. Connect it to all relevant hidden units in the literal layer with weights σ . Connect it to the rest of the hidden units in the literal layer with weights $+\beta$ or $-\beta$ with equal probabilities. Set the bias to $-\sigma^*(2n-1)/2$, where n stands for the number of relevant hidden units in the literal layer.
- 8 For each class, create an output unit and connect it to the relevant hidden units in the conjunction layer with weights σ . Connect it to the rest of the hidden units in the conjunction layer with weights $+\beta$ or $-\beta$ with equal probabilities. Set the bias to $-\sigma^*1/2$.

End

σ and β may be treated as parameters to this technique. The values for these were fixed by cross validating on an artificially created dataset. Values for momentum and learning rate were fixed similarly to 0.1 and 0.3. These values were subsequently frozen across datasets.

3 Example

Figure 1 depicts a typical decision tree that might be generated by running C4.5 on a training dataset comprising of examples with two attributes and two output classes. The classifier could be expressed equivalently, as the following rules in disjunctive normal form.

$$(X < 2.5) \vee ((X \geq 2.5) \wedge (Y < 1.3)) \Rightarrow \text{Class} : 1. \quad (3)$$

$$(X \geq 2.5) \wedge (Y \geq 1.3) \Rightarrow \text{Class} : 2. \quad (4)$$

The technique creates four nodes in the literal layer that correspond to $(X < 2.5)$, $(X \geq 2.5)$, $(Y < 1.3)$, and $(Y \geq 1.3)$. The literal $(X < 2.5)$ is characterized by the hidden node that has a bias of 12.5 and a weight of -5.0 on the input edge from X . The other nodes in the literal layer represent, in order, $(X \geq 2.5)$, $(Y < 1.3)$, and $(Y \geq 1.3)$.

Three nodes are generated in the conjunction layer; one each for $(X < 2.5)$, $(X \geq 2.5) \wedge (Y < 1.3)$ and $(X \geq 2.5) \wedge (Y \geq 1.3)$. For each such node, the weights on the edges from relevant nodes in the literal layer are set

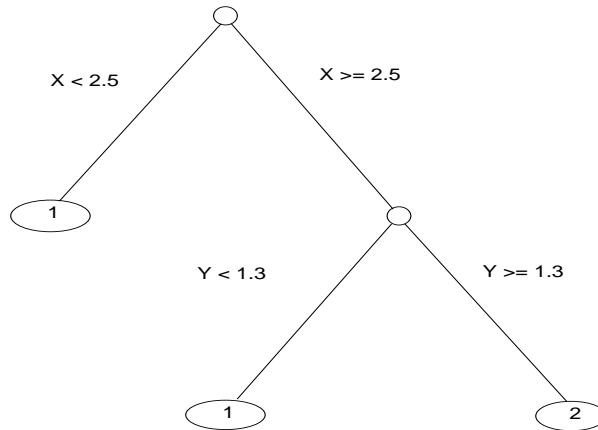


Figure 1: A typical decision tree.

to 5.0, and the bias of the node is set to $-(5n - 2.5)$, where n denotes the number of literals in the disjunct.

Two nodes are created in the output (disjunction) layer, one each for the two classes. The execution of this step is almost identical to the previous step. The only difference is that in this case, the bias of the nodes are set uniformly to -2.5 . A final step connects each node in a layer to nodes in the previous and the next layer that remain disconnected from it after the mentioned procedure. The weights on these edges are set to $+0.025$ and -0.025 with equal probabilities. Figure 2 depicts the skeletal network generated prior to the final step. Evidently, any classifier that can be reduced to a propositional ruleset of the form shown may be converted into an equivalent neural network using the given method.

4 Experiments and Results

The new technique was compared to C4.5 and two separate instances of **Bp**. The first network consisted of a single hidden layer with as many hidden units as were created by the new technique (strawman:1), and the second

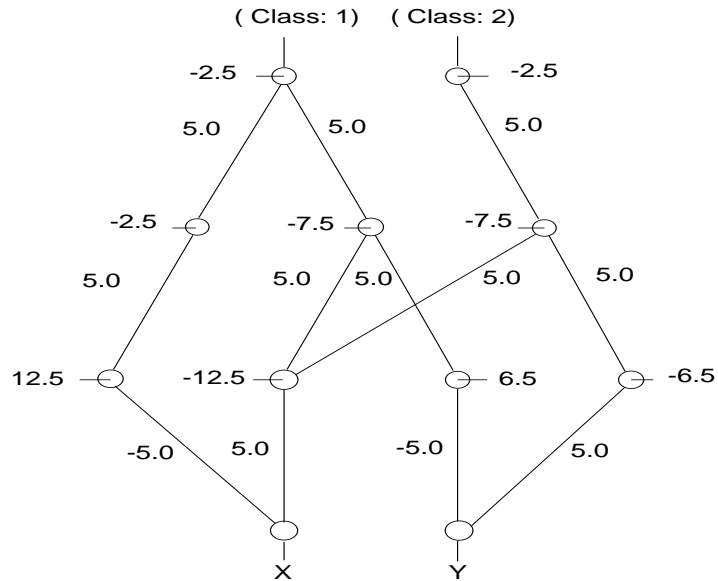


Figure 2: The skeletal neural network.

network employed the same topology as the network created by the technique (strawman:2). In both cases, weights on the edges were initialized randomly.

Experiments were conducted on three separate benchmark datasets. The first dataset was Iris, the second was a thyroid gland dataset donated to the UCI machine learning repository by Stefan Aeberhard and the third was a glass identification database donated to the UCI repository by Vina Spiehler.

In all experiments a 5-fold cross validation was performed with each training conducted thrice. **Bp** remained in a permutation mode of training through all the experiments. The values of momentum and learning rate for the network created by the technique were fixed at 0.1 and 0.3 respectively, across datasets. However, for the other networks, good values for momentum and learning rates were chosen. The number of epochs where the networks achieved their best predictive accuracies were also noted.

The results are presented graphically in figures 3 through 5, and in a tabular form in figure 6.

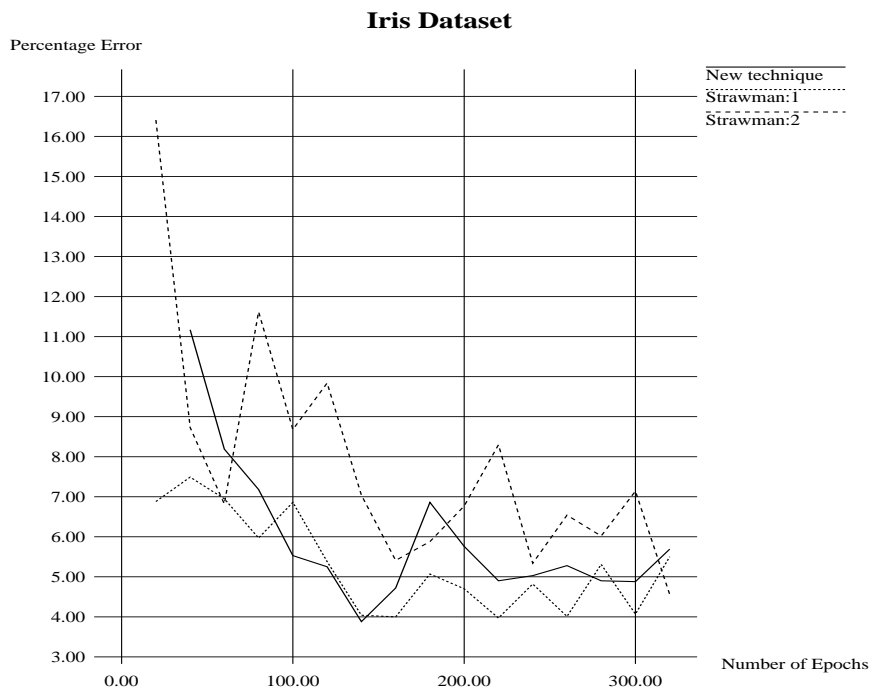


Figure 3: Results from the iris dataset.

4.1 Iris Dataset

When run on the Iris dataset, C4.5 generated a tree with an expected error rate of 8.5%. All three networks that were trained on the same dataset, however, surpassed C4.5 in accuracy for classifying novel examples. The new technique when trained to its optimal state, generated the best classifier with an expected error rate of 3.88%. Moreover, it reached its optimal state quicker than the other networks. Whereas strawman:1 and strawman:2 achieved their optimal states at 220 and 320 epochs respectively, the technique arrived at its optimal state at 140 epochs. Finally, the network with the same number of hidden units as used by the technique (strawman:1) performed better than the network with the same topology as the technique (strawman:2).

4.2 Thyroid Gland Dataset

In this case, the new technique performed distinctly better than C4.5 and the two other networks. The technique created a classifier with an expected error rate of 3.68% that surpassed, by far, the 4.37%, 5.08%, and 8.2% error rate

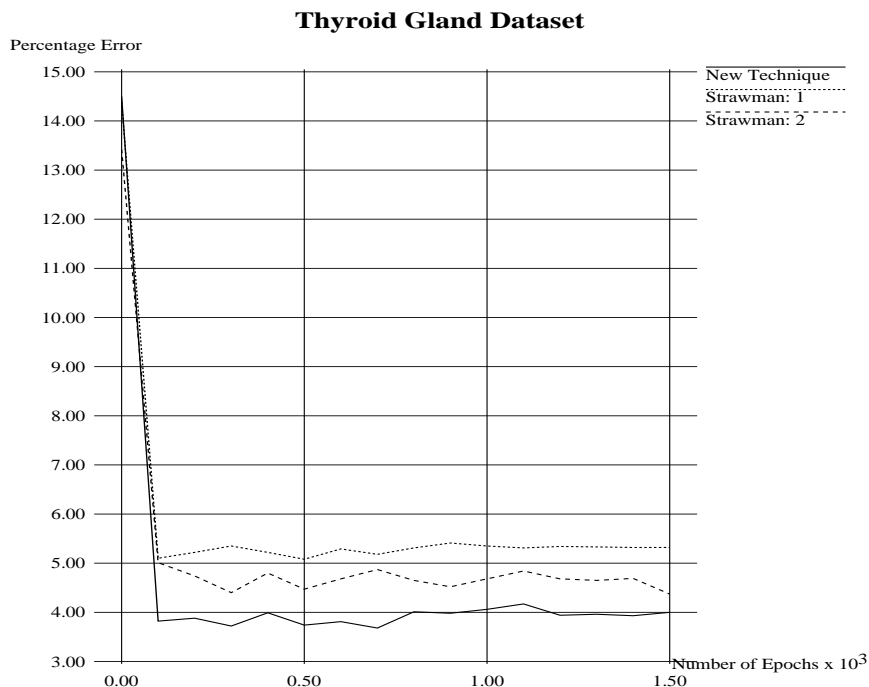


Figure 4: Results from the thyroid gland dataset.

classifiers created by strawman:2, strawman:1, and C4.5 respectively. Moreover, strawman:2 outperformed strawman:1 in this case unlike the previous instance.

4.3 Glass Identification Database

The Glass identification database has been known to be a dataset wherein C4.5 typically performs better than **Bp**. Results from experiments on this dataset, however, demonstrated that the new technique achieved a predictive accuracy that surpassed C4.5 and both the strawmen. This result could be attributed to the network’s initialization through C4.5.

One must note that the new technique performed better than the rival methods in spite of the basis for comparison being biased against it. Whereas the learning rate and momentum for the technique were frozen across datasets, good values for the parameters were chosen for the other methods.

Actual numerical values that were obtained through the experiments are recorded in a tabular form in figure 6. Results indicate that the complete

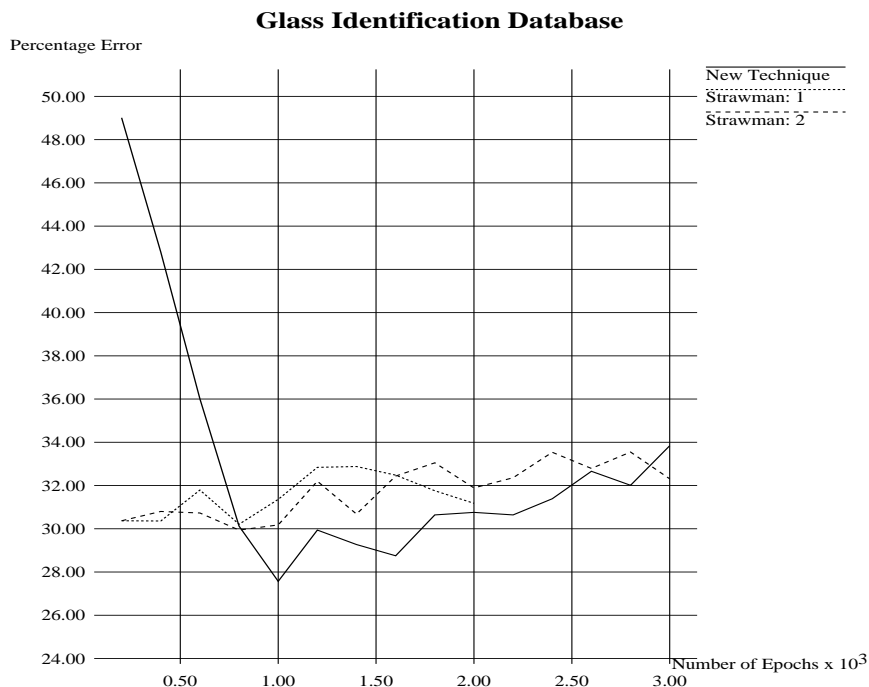


Figure 5: Results from the glass identification dataset.

process of finding the optimal values for momentum and learning rate through repeated training and cross validation can be discarded by applying, instead, the fixed parameter version of the new technique. Results also indicate that the efficacy of the technique can not be attributed to either the topology or the initial weights in isolation. It happens to be the consequence of the mutually supportive effects of both.

5 Future Research

In spite of the fact that the results from the experiments were quite promising, the new technique has been found to suffer from certain drawbacks. In all the experimented datasets, the decision trees generated by C4.5 were fairly small; trees that had fewer than 25 nodes. It is however, not very difficult to conceive of a dataset where any decision tree algorithm performs poorly and creates a complex tree with over a 100 internal nodes. Since, in the new technique the number of hidden units increase linearly with the number of internal nodes in the decision tree, it would not be at all surprising if the technique created a very large network topology for certain datasets.

	IRIS	THYROID GLAND DATASET	GLASS IDENTIFICATION DATABASE
C4.5	Expected error = 8.5%	Expected error= 8.2%	Expected error= 28.3%
Bp with same no: of hidden units	Best lrate = 0.7 Best momentum=0.5 Expected error= 3.97% after 220 epochs	Best lrate = 1.0 Best momentum=0.7 Expected error= 5.08% after 600 epochs	Best lrate = 0.8 Best momentum =0.6 Expected error= 30.21% after 800 epochs
Bp with same network topology	Best lrate = 0.7 Best momentum=0.5 Expected error= 4.56% after 320 epochs	Best lrate = 1.0 Best momentum= 0.7 Expected error= 4.37% after 1500 epochs	Best lrate = 0.8 Best momentum=0.6 Expected error= 29.94% after 800 epochs
New technique (Fixed lrate and momentum of 0.3 and 0.1 respectively)	Expected error= 3.88% after 140 epochs	Expected error= 3.68% after 700 epochs	Expected error= 27.57% after 1000 epochs

Figure 6: Complete tabulated results.

Applying the technique to multi-variate decision trees provides only a partial solution to this problem. Even though the hypothesis space for multi-variate decision trees is richer than that for uni-variate trees, it nevertheless remains weaker than the hypothesis space for a network with at least one hidden layer. This is to say that the initial network generated by any decision tree approach has many more nodes than is actually required to classify the dataset. The only solution to this problem lies in pruning the decision tree. Methods for pruning thus need to be investigated in this context.

6 Conclusions

To summarize, a host of drawbacks have been identified with the **Bp** algorithm. First, it is extremely slow to train. Second and more important, there

does not exist any formal technique that can decide beforehand the network topology and parameter settings that would be optimal for the network to classify a given dataset. This paper presents a novel technique that takes a decision tree classifier and transforms it into a network topology that classifies data almost as well as the tree. This process not only provides a good classifier to start with, but also chooses a suitable topology for the network. A short period of training then creates a classifier that is more accurate than **Bp**, and at the same time, runs faster. Finally, experimental comparisons against C4.5 and **Bp** demonstrate the efficacy of this new technique.

Acknowledgments

I would like to thank Haym Hirsh for stimulating discussions on this topic. I would also like to thank Jyoti Parmar for her critiques on an earlier draft of this paper.

References

- [1] Breiman, L., Friedman, J., Olshen, R., and Stone, C. (1984) Classification and Regression Trees. Wadsworth, Monterrey, Ca.
- [2] Brent, R., P. (1991) Fast Training Algorithms for Neural Nets. IEEE Transactions on Neural Networks, Volume 2.
- [3] Cios, K., J., and Liu, N. (1992) A Machine Learning Method for Generation of a Neural Network Architecture: A Continuous ID3 Algorithm. IEEE Transactions on Neural Networks, Volume 3.
- [4] Mooney, R., J., Shavlik, J., W., Towell, G., G., and Gove A. (1989) An Experimental Comparison of Symbolic and Connectionist Learning Algorithms. Proceedings of the 11th IJCAI.
- [5] Quinlan, J., R. (1986) Induction of Decision Trees. Machine Learning, Volume 1.
- [6] Rosenblatt, F. (1958) The Perceptron: A Probabilistic Model for Information Storage and Organization in the Brain. Psychological Review, Volume 65.

- [7] Rumelhart, D., E., Hinton, G., E., and Williams, R., J., (1986) Learning Internal Representations by Error Propagation. Parallel Distributed Processing Volume 1. Editors D. E. Rumelhart, and J. L. McClelland, MIT Press.
- [8] Sethi, I., K. (1990) Entropy Nets: From Decision Trees to Neural Networks. Proceedings of IEEE, Volume 78.
- [9] Towell, G., G., and Shavlik, J., W. (1993) Extracting Refined Rules from Knowledge-Based Neural Networks. Machine Learning, Volume 13.
- [10] Utgoff, P., E., (1989) Perceptron Trees: A case study in Hybrid Concept Representations. Connection Science, Volume 1.
- [11] Weiss, S., and Kapouleas, I. (1989) An Empirical Comparison of Pattern Recognition, Neural Nets, and Machine Learning Classification Methods. Proceedings of the 11th IJCAI.