

A Framework for Combining Symbolic and Neural Learning

Jude W. Shavlik

Computer Sciences Department
University of Wisconsin - Madison

shavlik@cs.wisc.edu

ABSTRACT

This article describes an approach to combining symbolic and connectionist approaches to machine learning. A three-stage framework is presented and the research of several groups is reviewed with respect to this framework. The first stage involves the insertion of symbolic knowledge into neural networks, the second addresses the refinement of this prior knowledge in its neural representation, while the third concerns the extraction of the refined symbolic knowledge. Experimental results and open research issues are discussed.

Keywords: knowledge-based neural networks
theory refinement
use of prior knowledge
rule extraction from neural networks
the KBANN algorithm
the NofM algorithm

A shorter version of this paper will appear in *Machine Learning*.

A Framework for Combining Symbolic and Neural Learning

Jude W. Shavlik

Computer Sciences Department
University of Wisconsin - Madison

Introduction

The last ten or so years have produced an explosion in the amount of research on machine learning. This rapid growth has occurred, largely independently, in both the symbolic and connectionist (neural network) machine learning communities. Fortunately, over the last few years these two communities have become less separate, and there has been an increasing amount of research that can be considered a hybrid of the two approaches. This paper reviews some of the research that combines the symbolic and neural network approaches to artificial intelligence and presents a framework for combining the two paradigms.

We will not attempt to define precisely the essential differences between the symbolic and connectionist approaches, as that would lead to a lengthy debate far beyond the scope of this article. If some distinction is needed, we can make the coarse approximation that symbolic approaches focus on producing discrete combinations of features, while neural approaches adjust continuous, non-linear weightings of their inputs. However, we will assume that understanding the fundamental differences between the two paradigms is a future research issue, and we will focus on research that incorporates what traditionally might be considered aspects of both camps.

There are a large number of ways to combine symbolic and connectionist AI. For example, Utgoff (1988) developed an algorithm that closely integrates decision trees and perceptrons. One could also have a loosely-coupled hybrid system in which "high" level decisions are made symbolically, while "low" level one are made by neural networks (e.g., Gallant, 1988; Pomerleau, Gowdy, & Thorpe, 1991). Recent special issues of journals (Hendler, 1989; Hinton, 1990) present additional approaches. However, rather than attempting a comprehensive review of all the symbolic/connectionist hybrid methods explored, we will focus on the framework that Figure 1 illustrates.

In this framework, the learner first *inserts* symbolic information of some sort into a neural network; it is becoming increasingly clear that a learner must make effective use of prior knowledge in order to perform well (Geman, Bienenstock, & Doursat, 1992). Once in a neural representation, it uses training examples to *refine* the initial knowledge. Finally, it *extracts* symbolic information from the trained network. The research of several groups fits nicely into this framework, and promising results have been achieved. The remainder of this paper discusses some of this research and points out open issues in each of the three phases. But before continuing, it should be noted that these three steps are somewhat independent and researchers have studied various combinations of them.

The remainder of this article is organized around four questions. We first consider why one should use neural networks for symbol-oriented learning tasks. We then review research that addresses questions about each of Figure 1's three phases: insertion, refinement, and extraction of symbolic information.

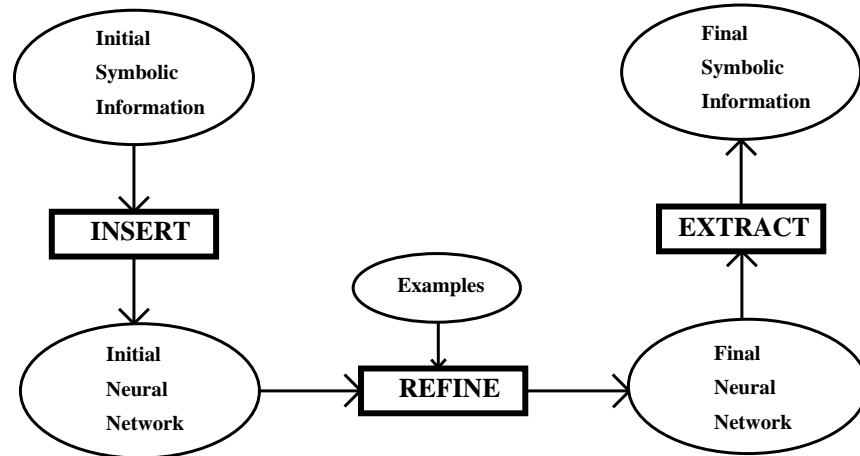


Figure 1. A framework for combining symbolic and neural learning.

Why use Neural Networks for Symbol-Oriented Learning Tasks?

Should one avoid using connectionist methods to learn tasks that inherently deal with symbols? Are not neural networks primarily applicable to "low-level", perceptual tasks? We will argue in this section that the answer to these related questions is "no."

Over the last few years, starting with three papers published simultaneously at IJCAI-89 (Fisher & McKusick, 1989; Mooney et al., 1989; Weiss & Kapouleas, 1989) and followed by other studies (e.g., Atlas, 1990; Dietterich, Hild, & Bakiri, 1990), several groups have empirically compared symbolic learning algorithms, such as Quinlan's (1986) ID3 decision-tree algorithm, to connectionist approaches, such as Rumelhart, Hinton, and Williams' (1986) backpropagation method for training neural networks. These studies did not produce consistent results, but their coarse summary is that trained neural networks have at least comparable accuracies to induced decision trees on tasks that can be considered symbol oriented. Hence, it appears worthwhile to investigate using neural learning methods to produce and refine symbolic information.

In addition, neural network approaches have proven successful on a wide range of "real world" tasks, such as speech understanding (Lippmann, 1989), handwritten-character recognition (Le Cun et al., 1989), control of dynamic systems (Jordan & Rumelhart, 1992), gene finding (Uberbacher & Mural, 1991), and language learning (Touretzky, 1991). These experiments strongly suggest that connectionist learning is a powerful approach, and the use of neural networks with symbolic knowledge merits exploration.

Finally, it is important to note that there are connectionist architectures beyond the simple, feed-forward, single-hidden-layer neural networks. In particular, recurrent networks (Elman, 1990; Jordan, 1986), with their feedback loops and "memory", are especially appealing for application to symbolic tasks that have a sequential nature.

How can We Get Symbolic Information into Neural Networks?

Assuming one is convinced of the merit of Figure 1's framework, techniques for inserting symbolic information into a neural network are needed. One can think of this preexisting information as prior knowledge about the task at hand, and the question is: how can neural networks effectively use these "hints" (Abu-Mostafa, 1990)?

One answer, the KBANN approach (Towell, Shavlik, & Noordewier, 1990; Towell, 1992), creates *knowledge-based artificial neural networks* by producing neural networks whose topological structure matches the dependency structure of the rules in an approximately-correct "domain theory" (a collection of inference rules about the current task). Table 1 shows the correspondences between a domain theory and a neural network, and Figure 2 contains a simple example of the K approach to mapping a domain theory into a neural networks.

KBANN has been applied to successfully refining domain theories for real-world problems such as gene finding (Towell et al., 1990), protein folding (Maclin & Shavlik, in press), and the control of a simple chemical plant (Scott, Shavlik, & Ray, 1992). The appendix to this paper presents the application of KBANN to a problem from the Human Genome Project. The appendix' example shows a more complicated mapping than that in Figure 2 and also reports how well the KBANN approach generalizes to examples not seen during training.

Various groups have found that knowledge-based neural networks train faster than do "standard" neural networks (Berenji, 1991; Oliver & Schneider, 1988; Omlin & Giles, 1992; Shavlik & Towell, 1989), presumably because the initial information is used to choose a good starting point for the network. More importantly, though, experiments have shown that knowledge-based networks generalize better to future examples than do standard networks, as well as several other methods for inductive learning and theory refinement (Omlin & Giles, 1992; Maclin & Shavlik, in press; McMillan et al., 1992; Roscheisen, Hofmann, & Tresp, 1992; Scott et al., 1992; Towell, 1992; Towell et al., 1990; Tresp, Hollatz, & Ahmad, 1993). One can attribute this improved generalization to two aspects of the insertion process. The domain theory produces a useful inductive bias by (1) focusing attention on relevant input features and (2) indicating useful intermediate conclusions (which suggest a good network topology).

Figure 3 schematically illustrates the general performance of the KBANN system. It plots generalization performance (i.e., testset accuracy) as a function of the number of examples

Table 1. Correspondences between a domain theory and a neural network.

Domain Theory	Neural Network
final conclusion	output units
intermediate conclusions	hidden units
supporting facts	input units
antecedents of a rule	highly-weighted links

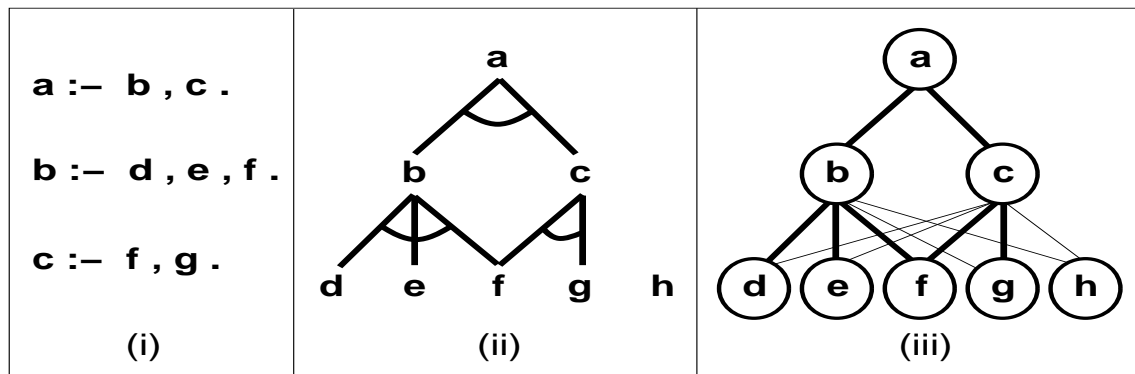


Figure 2. A sample application of the KBANN rule-insertion algorithm. Frame (i) contains a simple domain theory, while frame (ii) shows the dependency structure of these rules. The third frame shows the network KBANN creates. The thick lines in (iii) correspond to the dependencies in the rules; KBANN sets the weights on these links in such a manner that nodes are highly active only when the domain theory supports the corresponding deduction. Thin lines in frame (iii) represent zero-weighted links that KBANN adds to the network to allow refinement of the domain theory during neural training.

available for training. Qualitatively similar results were obtained on four different domains (see Towell et al., 1990, Noordewier, Towell, and Shavlik, 1991, Scott et al., 1992 and Maclin & Shavlik, in press, for the specific generalization curves on the four domains). In each case, KBANN was given an imperfect domain theory. By using this domain-specific information, KBANN produced a good starting network, which, for small numbers of training examples, resulted in statistically-significantly better testset accuracy than that obtained by a standard artificial neural network (ANN) initialized with small, random weights.

However, asymptotically the testset accuracy of the two learning systems converged on those experimental domains where we could collect enough training examples. This asymptotic convergence suggests another view of the value of a domain theory: the initial knowledge is "worth" some number of training examples. In many domains, collecting a large number of examples is impossible or at least very costly (e.g., protein folding), so being able to utilize alternate sources of information can prove quite valuable.

Towell (1992) has shown that KBANN's knowledge-based networks better refine a domain theory than do purely symbolic theory-refinement systems. This holds even when one compares the rules extracted from the trained network to the refined rules produced by the symbolic theory-refinement systems; these results provide a justification for the complex representational shifts in Figure 1's framework. One can convert the rules that KBANN extracts to disjunctions of conjunctive rules (usually with a great increase in the number of rules), so that the two approaches are searching the same hypothesis space. While Towell's empirical results may well

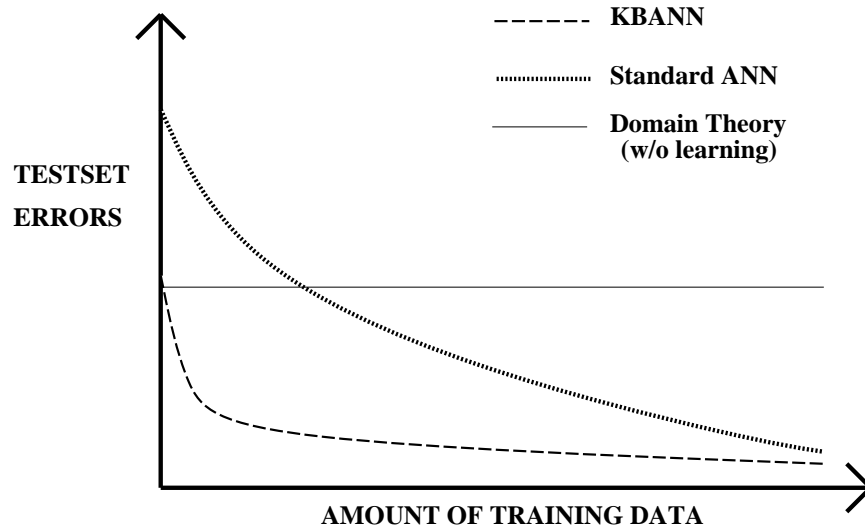


Figure 3. Generalization to new examples as a function of the number of training examples.

be problem-specific, a broader conclusion is that searching the continuous weight space of neural networks is better on "real-world" problems than searching the combinatorial space of discrete rules - complex concepts in one representation may be much simpler in the other. A deeper understanding of the relative merits of the symbolic/connectionist and the purely symbolic approaches to theory refinement is an important open research issue.

In addition to the simple, propositional rules shown in Figure 2 and used in much of the early KBANN work, researchers have produced techniques for mapping several other forms of prior knowledge into networks. Fu (1989) and Mahoney and Mooney (1993) map rules containing certainty factors. Berenji (1991) and Masuoka et al. (1990) map fuzzy-logic rules, while McMillan, Mozer, and Smolensky (1992) use gating networks (Jacobs et al., 1991) to map production rules. Scott et al. (1992) and Roscheisen et al. (1992) map mathematical equations, demonstrating that the KBANN approach does not require logic-oriented domain theories. Finally, several groups have mapped (generalized) finite-state grammars into recurrent neural networks (Fransconi et al., 1991; Maclin & Shavlik, in press; Omlin & Giles, 1992; Scott et al., 1992). Generalized finite-state grammars are particularly interesting to the theory-refinement community, as one can view them as state-dependent domain theories, a richer type of domain theory than is usually studied in this subfield of machine learning. These approaches differ from KBANN to various degrees, but the essential idea is the same: use prior knowledge to decide how to initialize a neural network.

Most of the domain theories studied in the theory-refinement literature are propositional and address simple classification tasks. As mentioned above, several groups have been studying the use of *finite-state automata* as a means of expressing approximately-correct domain knowledge. "Finite-state" domain theories allow a system to have a memory, which means it

Combining Symbolic and Neural Learning

can base its decisions not only on the current input but on its summary of recent inputs as well. Figure 4 shows the architecture of a connectionist approach to refining these state-dependent domain theories. Using recurrent neural networks, the system bases its output on both the current input and its internal state (e.g., a "world model" in planning problems). In addition to calculating the current output, the system must also select its next state. A state-dependent domain theory is used to initially configure the network, in a manner analogous to that shown in Figure 2.

There are several open questions regarding the knowledge-insertion process. We would like to know what other types of prior knowledge can be inserted into networks. For example, methods are lacking for inserting first-order theories. The last few years have seen much progress in inductive logic programming (Muggleton, 1992; Quinlan, 1990), and it would be useful to see if (and how well) neural networks can refine rules containing variables. To do so, one needs to devise methods for dealing with unbounded symbolic structures in neural networks (whose size is usually fixed following training). Recurrent networks provide one method of dealing with unbounded structures, and Pollack's (1990) recursive auto-associative memories provide another. Also relevant is research on teaching networks to recognize context-free grammars by having them learn how to use a stack (Das et al., 1993; Giles et al., 1990; Mozer & Das, 1993). Unbounded structures, such as stacks, can be handled in fixed-size networks by somehow altering resolution (in some sense) so that the "product" of the information being stored and its resolution equals a constant.

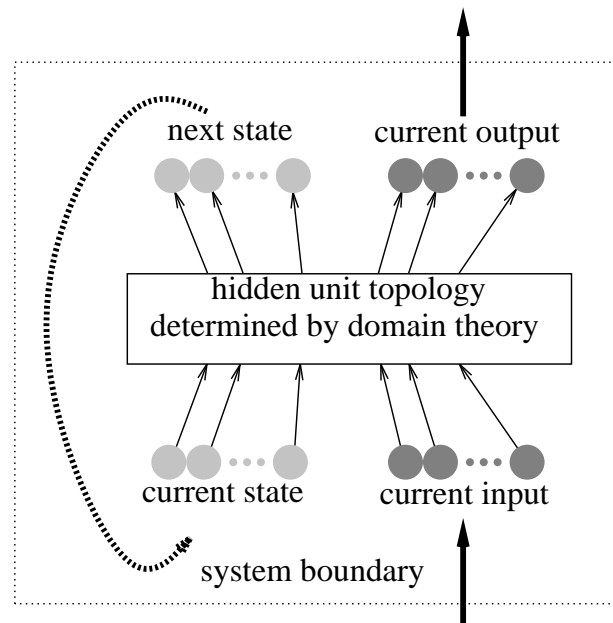


Figure 4. Refining "finite-state" domain theories.

Combining Symbolic and Neural Learning

Towell (1992) has shown that knowledge-based networks are good at deleting irrelevant information in approximately-correct domain theories, but do not handle "impoverished" domain theories as well. Figure 5 shows the qualitative nature of his experiments, which demonstrate that domain theories need only be approximately correct to prove beneficial. Towell added spurious antecedents to the rules of an existing domain theory, and found that until there was a 50% expansion, the corrupted domain theory still lead to better generalization than a standard, one hidden-layer network (which does not use the domain theory). However, he could only delete 25% of the antecedents before the domain theory lead to worse testset accuracy. In other words, the KBANN approach is better at discarding erroneous information than discovering missing knowledge. Hence, another open issue in knowledge-based neural networks is how to deal with domain theories that are incomplete. We will return to this topic in the next section.

Converting symbolic information to a neural-network representation, followed by connectionist learning, has been shown useful by several research groups. This leaves us with the central question about the insertion phase:

How can we convert symbolic knowledge and learning tasks so that powerful numeric optimization search methods are applicable?

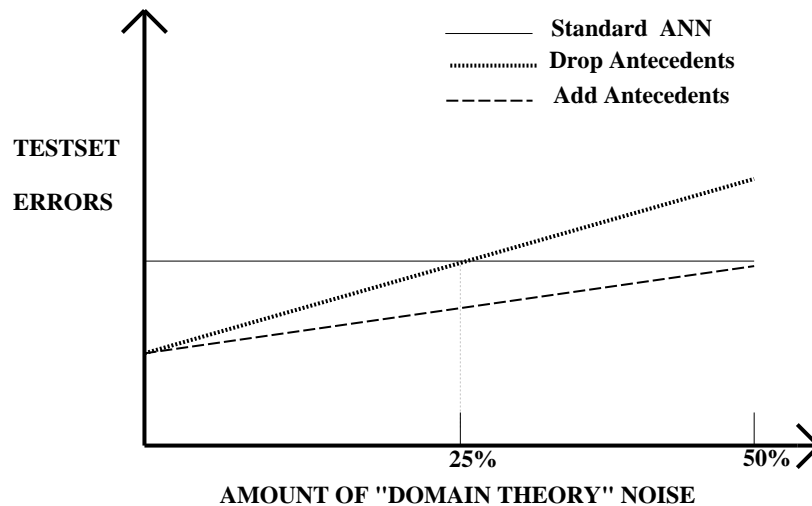


Figure 5. The effect on generalization of corrupting the domain theory.

How can Network Refinement be Guided by Symbolic Knowledge?

Once prior knowledge is inserted into a network, it has to be refined and enhanced. A simple way of doing this is to run backpropagation, or some other standard connectionist training procedure, on the training examples. However, there are two ways to use symbolic information to improve training: (1) one could use symbolic learning methods and ideas to focus the adjustment of the network, both its weights and topology, and (2) one might alter backpropagation to better match the symbolic nature of a given problem. In this section we discuss both of these approaches.

One might ask, where is the symbolic learning in the approaches presented so far? One answer is that domain-theory refinement, which is what knowledge-based networks do, addresses the incorrect-theory problem of explanation-based learning. In fact, this perspective was the initial motivation for the KBANN research (Shavlik & Towell, 1989). (Recently, Mitchell and Thrun (1993) proposed an explanation-based, though non-symbolic, method for training neural networks for reinforcement learning tasks.)

But what about performing symbolic inductive learning in conjunction with neural learning? As mentioned above, Utgoff's (1988) perceptron trees are one method for doing so, but his algorithm is not applicable to the refinement of prior knowledge. Recall that in knowledge-based networks, input features fall into two classes: those that are mentioned in the domain theory and those that are not. Since the domain theory can be imperfect, one cannot ignore the unmentioned input features; they are typically connected to other units with low-weighted links. Towell and Shavlik (1992a) proposed a technique that uses symbolic inductive learning to identify good input features, which are then weighted more heavily. They found this preprocessing of the network led to better generalization.

Also, as mentioned, a domain theory may be missing a number of rules. Hence, it will be mapped into a network that is too small. In order to learn these missing rules, additional nodes will have to be added to the network during training. Opitz and Shavlik (1992) developed an algorithm that interprets networks symbolically to decide where to add new nodes.

There have been several changes to standard connectionist learning motivated by symbolic problems. Rather than minimizing mean-squared error, the cross-entropy error function (Hinton, 1989) is a better choice for knowledge-based networks (see Towell (1992) for an explanation). Refining rules with certainty factors requires the use of a different activation function for nodes (Fu, 1989; Mahoney & Mooney, 1993). One may wish to constrain weight changes to maintain the symbolic interpretation of the network (McMillan et al., 1992). Finally, networks often decay their weights toward zero during training (Hinton, 1986). Weights in knowledge-based networks should decay toward their initial values (Hinton, personal communication; Tresp et al., 1993), thereby encouraging the network to preserve the knowledge in the initial domain theory.

There are several open questions regarding the use of symbolic information to aid the refinement step. How can one detect that extra nodes are needed to generalize well, and where are the best places to add them? Folk wisdom says that backpropagation does not work well in networks with many layers of hidden units, because the error signal becomes too diffuse. Can one use symbolic information to focus the back-propagated error signal, especially in deep networks? Deep networks often occur when basing the network topology on the dependency structure of a rule base, so this problem is exacerbated in knowledge-based networks. Finally, do we need to prevent distributed representations (Hinton, 1986) from evolving during training? Since hidden units in knowledge-based networks initially have a symbolic meaning, it seems that

distributed representations are undesirable; however, there could be some way to take advantage of distributed representations.

In summary, the central question about the refinement phase is:

How can symbolic knowledge about
the task at hand guide network refinement?

How can We Extract Symbolic Knowledge from Trained Neural Networks?

The third phase of Figure 1's framework involves extracting symbolic information (e.g. rules) from a trained network, which need not originally be knowledge-based. Why is this important? Rule extraction can help one understand what the "black box" network has learned. If the network produced a scientifically-interesting discovery, it would be nice if this were made explicit. Also, one may wish that a trained system would produce explanations of its future decisions. Finally, one may want to manipulate the results of learning in another system, such as a planner.

Several people have developed methods for extracting rules from standard networks. Gallant (1988), Saito and Nakano (1988), and Fu (1991) proposed algorithms that consider various ways that a node's weighted input can exceed its threshold, and convert each of these situations into a rule. However, these approaches can require an exponential number of rules (in terms of the number of network weights) to re-represent a node. Towell and Shavlik (1992b) developed a method that produces about one "N out of M" rule for each node. They found their algorithm extracted comprehensible rules while maintaining the accuracy of the trained network. However, their approach only works well on knowledge-based networks, as it requires that weights cluster into a few groups; the "soft-weight sharing" technique of Nowlan and Hinton (1992) may improve the performance of Towell and Shavlik's algorithm on standard networks. Finally, McMillan et al. (1992) simply project trained nodes to the closest valid rule, while Hayashi (1991) extracts a small number of fuzzy-logic rules from a trained network.

Figure 6 provides a sketch of Towell and Shavlik's NofM algorithm. This algorithm extracts rules from trained knowledge-based networks, and it assumes that the nodes in a trained network always have values near zero or one (this "Boolean constraint" is easily achieved by steepening the sigmodal activation function during the final stages of training). The first step - see the upper, middle panel in Figure 6 - clusters the incoming weights of a unit, using a standard clustering algorithm (Hartigan, 1975). The weights on the links in each cluster are then replaced by the average of the cluster's weights, which regularizes the network. Next, the algorithm analyzes each cluster to determine if any clusters are irrelevant. In Figure 6, the cluster whose average weight value is 1.1 can be discarded, because, due to the Boolean constraint stated above, the activation of unit *Z* will never be qualitatively effected by the settings of nodes *B*, *D*, *E*, and *G*. A simple case analysis demonstrates this. If zero or one of $\{A, C, F\}$ are active, the maximum weighted input into node *Z* is 10.5. While if two or three of $\{A, C, F\}$ are active, the minimum weighted input into node *Z* is 12.2. In neither case does the specific activations of $\{B, D, E, G\}$ have an impact. The final step is to rewrite the regularized and simplified node as a rule, as shown in the bottom panel of the figure. The NofM algorithm does not always produce rules as simple and clean as the one in Figure 6 - see Towell and Shavlik (1992b) for a real-world example - but this figure does illustrate the essential aspects of the NofM rule-extraction technique.

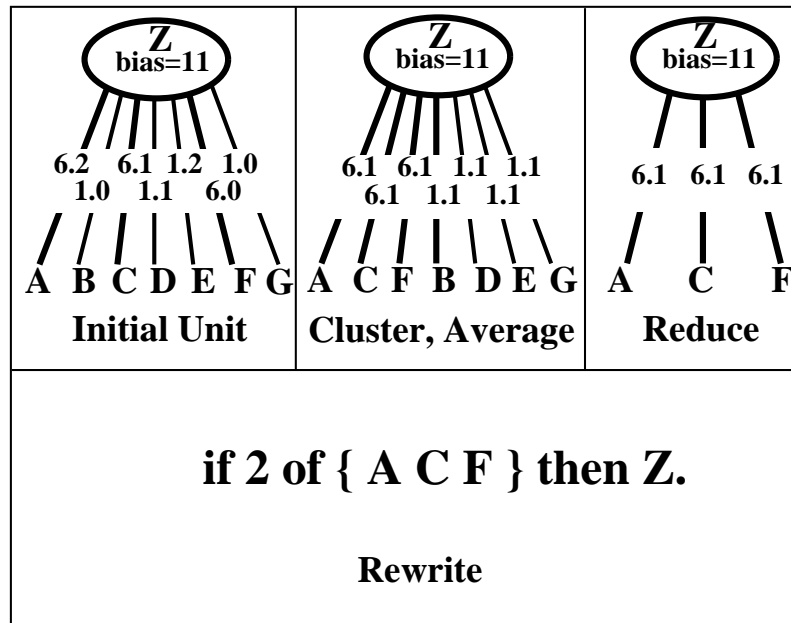


Figure 6. A sketch of the NofM rule-extraction algorithm.

The above methods analyze the weights going into nodes. Cleermans, Servan-Schreiber, and McClelland (1989) and Giles et al. (1992) have a different perspective. They investigate extracting finite-state automata from recurrent networks, and their methods focus on the activation patterns of the hidden units. Their approaches assume that these patterns represent some sort of internal state. The extraction algorithms cluster these patterns and view each cluster as a state in their automaton. The next step runs the training examples through the trained network to obtain the state transitions, after which traditional algorithms minimize the automaton.

A major question with rule extraction is: how does one measure comprehensibility? An extraction algorithm must produce reasonably comprehensible rules, but without a good measure it is hard to compare alternative approaches. A second open issue relates to the refinement phase: how should this task be altered in support of rule extraction? Possibly the network can be constrained to always lie in the "comprehensible" portion of weight space, whatever that might be. Related to this, the hidden units in knowledge-based networks generally have symbolic names attached to them, and if one is going to use these labels for the extracted rules, one needs to ensure that the symbol-node correspondence is not altered during training. This is one reason why the formation of distributed representations during training can be harmful. Finally, conceptually clustering hidden-unit activations is a promising area for symbolic machine learning. Finding and describing clusters can provide insight into the distinctions made by the network. For example, Sejnowski and Rosenberg (1987) manually analyzed clusters developed on the NETtalk task, with some success.

Combining Symbolic and Neural Learning

To wrap up this section, the central question about rule extraction is:

How can we extract a small and comprehensible symbolic version of a trained network without losing (much) accuracy?

Conclusion

Connectionist machine learning has proven to be a fruitful approach, and it makes sense to investigate systems that combine the strengths of the symbolic and connectionist approaches to AI. Over the past few years, researchers have successfully developed a number of such systems. This article summarizes one view of this endeavor, a framework that encompasses the approaches of several different research groups. This framework (see Figure 1) views the combination of symbolic and neural learning as a three-stage process: (1) the insertion of symbolic information into a neural network, thereby (partially) determining the topology and initial weight settings of a network, (2) the refinement of this network using a numeric optimization method such as backpropagation, possibly under the guidance of symbolic knowledge, and (3) the extraction of symbolic rules that accurately represent the knowledge contained in a trained network. These three components form an appealing, complete picture - approximately-correct symbolic information in, more-accurate symbolic information out - however, these three stages can be independently studied. In conclusion, the research summarized in this paper demonstrates that combining symbolic and connectionist methods is a promising approach to machine learning.

Acknowledgements

The material in this article is based on an invited talk presented at the 1992 International Machine Learning Conference held in Aberdeen, Scotland. I wish to thank Geoff Towell, Mick Noordewier, Rich Maclin, Gary Scott, Mark Craven, Dave Opitz, Derek Zahn, Charlie Squires, and Kevin Cherkauer - all members of the University of Wisconsin Machine Learning Research Group at one time or another and major contributors to the ideas presented in this chapter. Discussions with Ray Mooney, Tom Dietterich, and Geoff Hinton also substantially influenced this discussion. This work is partially supported by Office of Naval Research Grant N00014-90-J-1941, National Science Foundation Grant IRI-9002413, and Department of Energy Grant DE-FG02-91ER61129.

References

- Atlas, L., Cole, R., Connor, J., El-Sharkawi, M., Marks II, R., Muthusamy, Y., & Barnard, E. (1990). Performance comparisons between backpropagation networks and classification trees on three real-world applications. In *Advances in Neural Information Processing Systems (Vol. 2)*, D. Touretzky (ed.), San Mateo, CA: Morgan Kaufmann.
- Abu-Mostafa, Y. S. (1990). Learning from hints in neural networks. *Journal of Complexity*, 6, 192-198.
- Berenji, H. R. (1991). Refinement of approximate reasoning-based controllers by reinforcement learning. *Proceedings of the Eighth International Machine Learning Workshop* (pp. 475-479), Evanston, IL: Morgan Kaufmann.
- Cleermans, A., Servan-Schreiber, D., & McClelland, J. L. (1989). Finite state automata and simple recurrent networks. *Neural Computation*, 1, 372-381.

Combining Symbolic and Neural Learning

- Das, S., Giles, C. L., & G. Z. Sun (1993). Using hints to successfully learn context-free grammars with a neural network pushdown automaton. In *Advances in Neural Information Processing Systems (Vol. 5)*, S. Hanson, J. Cowans, & L. Giles (eds.), San Mateo, CA: Morgan Kaufmann.
- Dietterich, T. G., Hild, H., & Bakiri, G. (1990). A comparative study of ID3 and backpropagation for English text-to-speech mapping. *Proceedings of the Seventh International Conference on Machine Learning* (pp. 24-31), Austin, TX: Morgan Kaufmann.
- Elman, J. L. (1990). Finding structure in time. *Cognitive Science*, *14*, 179-211.
- Fisher, D. H. & McKusick, K. B. (1989). An empirical comparison of ID3 and back-propagation. *Proceedings of the Eleventh International Joint Conference on Artificial CAIntelligence* (pp. 788-793). Detroit: Morgan Kaufmann.
- Frasconi, P., Gori, M., Maggini, M., & Soda, G. (1991). A unified approach for integrating explicit knowledge and learning by example in recurrent networks. *Proceedings of the International Joint Conference on Neural Networks*, (pp. 811-816). Seattle: IEEE Press.
- Fu, L. M. (1989). Integration of neural heuristics into knowledge-based inference. *Connection Science*, *1*, 325-340.
- Fu, L. M. (1991). Rule learning by searching on adapted nets. *Proceedings of the Ninth National Conference on Artificial Intelligence* (pp. 590-595). Anaheim, CA: AAAI Press.
- Gallant, S. I. (1988). Connectionist expert systems. *Communications of the ACM*, *31*, 152-169.
- Gemen, S., Bienenstock, E. & Doursat, R. (1992). Neural networks and the bias/variance dilemma. *Neural Computation*, *4*, 1-58.
- Giles, C., Sun, G., Chen, H., Lee, Y., & Chen, D. (1990). Higher order recurrent networks and grammatical inference. In *Advances in Neural Information Processing Systems (Vol. 2)*, D. Touretzky (ed.), San Mateo, : Morgan Kaufmann.
- Giles, C., Miller, C., Chen, D., Chen, H., Sun, G., & Lee, Y. (1992). Learning and extracting finite state automata with second-order recurrent neural networks. *Neural Computation*, *4*, 393-405.
- Hartigan, J. A. (1975). *Clustering Algorithms*, New York: Wiley.
- Hawley, D. K. & McClure, W. R. (1983). Compilation and analysis of Escherichia Coli promoter DNA sequences. *Nucleic Acids Research*, *11*, 2237-2255.
- Hayashi, Y. (1991). A neural expert system with automated extraction of fuzzy if-then rules and its application to medical diagnosis. In *Advances in Neural Information Processing Systems (Vol. 3)*, R. Lippmann, J. Moody, & D. Touretzky (eds.), San Mateo, CA: Morgan Kaufmann.
- Hendler, J. A. (ed.) (1989). Special issue on hybrid systems (symbolic/connectionist). *Connection Science*, *1*.
- Hinton, G. E. (1986). Learning distributed representations of concepts. *Proceedings of the Eighth Annual Conference of the Cognitive Science Society* (pp. 1-12). Amherst, MA: Erlbaum.
- Hinton, G. E. (1989). Connectionist learning procedures *Artificial Intelligence*, *40*, 185-234.
- Hinton, G. E. (ed.) (1990). Special issue on connectionist symbol processing. *Artificial Intelligence*, *46*.
- Jacobs, R. A., Jordan, M. I., Nowlan, S. J., & Hinton, G. E. (1991). Adaptive mixtures of local experts. *Neural Computation*, *3*, 79-87.
- Jordan, M. I. (1986). Attractor dynamics and parallelism in a connectionist sequential machine. *Proceedings of the Eighth Annual Conference of the Cognitive Science Society*, (pp. 531-546), Amherst, MA: Erlbaum.
- Jordan, M. I. & Rumelhart, D. E. (1992). Forward models: Supervised learning with a distal teacher. *Cognitive Science*, *16*, 307-354.
- Le Cun, Y., Boser, B., Denker, J., Henderson, D., Howard, R., Hubbard, W., & Jackel, L. (1989). Backpropagation applied to handwritten zip code recognition. *Neural Computation*, *1*, 541-551.
- Lippmann, R. P. (1989). Review of neural networks for speech recognition. *Neural Computation*, *1*, 1-38.
- Maclin, R. & Shavlik, J. W. (in press). Using knowledge-based neural networks to improve algorithms: Refining the Chou-Fasman algorithm for protein folding. *Machine Learning*.
- Mahoney, J. J. & Mooney, R. J. (1993). Combining neural and symbolic learning to revise probabilistic rule bases. In *Advances in Neural Information Processing Systems (Vol. 5)*, S. Hanson, J. Cowans, & L. Giles (eds.), San Mateo, CA: Morgan Kaufmann.

Combining Symbolic and Neural Learning

- Masuoka, R., Watanabe, N., Kawamura, A., Owada, Y., & Asakawa, K. (1990). Neurofuzzy system - fuzzy inference using a structured neural network. *Proceedings of the International Conference on Fuzzy Logic & Neural Networks* (pp. 173-177), Iizuka, Japan.
- McMillan, C., Mozer, M. C., & Smolensky, P. (1992). Rule induction through integrated symbolic and subsymbolic processing. In *Advances in Neural Information Processing Systems (Vol. 4)*, J. Moody, S. Hanson, & R. Lippmann (eds.), San Mateo, CA: Morgan Kaufmann.
- Mitchell, T. M. & Thrun, S. (1993). Explanation-based neural network learning for robot control. In *Advances in Neural Information Processing Systems (Vol. 5)*, S. Hanson, J. Cowans, & L. Giles (eds.), San Mateo, CA: Morgan Kaufmann.
- Mooney, R., Shavlik, J., Towell, G., & Gove, A. (1989). An experimental comparison of symbolic and connectionist learning algorithms. *Proceedings of the Eleventh International Joint Conference on Artificial Intelligence* (pp. 775-780). Detroit: Morgan Kaufmann. (An extended version appeared in *Machine Learning*, 6, 111-143, 1991.)
- Mozer, M. C. & Das, S. (1993). A connectionist chunker that induces the structure of context-free languages. In *Advances in Neural Information Processing Systems (Vol. 5)*, S. Hanson, J. Cowans, & L. Giles (eds.), San Mateo, CA: Morgan Kaufmann.
- Muggleton, S. (1992). *Inductive Logic Programming*, London: Academic Press.
- Noordewier, M. O., Towell, G. G., & Shavlik, J. W. (1991). Training knowledge-based neural networks to recognize genes in DNA sequences. In *Advances in Neural Information Processing Systems (Vol. 3)*, R. Lippmann, J. Moody, & D. Touretzky (eds.), San Mateo, CA: Morgan Kaufmann.
- Nowlan, S. J. & Hinton, G. E. (1992). Simplifying neural networks by soft weight-sharing. In *Advances in Neural Information Processing Systems (Vol. 4)*, J. Moody, S. Hanson, & R. Lippmann (eds.), San Mateo, CA: Morgan Kaufmann.
- Oliver, W. L. & Schneider, W. (1988). Using rules and task division to augment connectionist learning. *Proceedings of the Tenth Annual Conference of the Cognitive Science Society* (pp. 55-61), Montreal: Erlbaum.
- Omlin, C. W. & Giles, C. L. (1992). Training second-order recurrent neural networks using hints. *Proceedings of the Ninth International Conference on Machine Learning* (pp. 361-366), Aberdeen, Scotland: Morgan Kaufmann.
- O'Neill, M. C. (1989). Escherichia Coli promoters. *Journal of Biological Chemistry*, 264, 5522-5530.
- Opitz, D. & Shavlik, J.W. (1992). *Heuristically expanding knowledge-based neural networks*. Technical Report, Madison, WI: University of Wisconsin, Computer Sciences Department.
- Pollack, J. (1990). Recursive distributed representations. *Artificial Intelligence*, 46, 77-105.
- Pomerleau, D.A., Gowdy, J., & Thorpe, C.E. (1991). Combining artificial neural networks and symbolic processing for autonomous robot guidance. In *Engineering Applications of Artificial Intelligence*, 4, 279-285.
- Quinlan, J. R. (1986). Induction of decision trees. *Machine Learning*, 1, 81-106.
- Quinlan, J. R. (1990). Learning logical definitions from relations, *Machine Learning*, 5, 239-266.
- Roscheisen, M., Hofmann, R. & Tresp, V. (1992). Neural control for rolling mills: Incorporating domain theories to overcome data deficiency. In *Advances in Neural Information Processing Systems (Vol. 4)*, J. Moody, S. Hanson, & R. Lippmann (eds.), San Mateo, CA: Morgan Kaufmann.
- Rumelhart, D. E., Hinton, G. E., & Williams, R. J. (1986). Learning internal representations by error propagation. In *Parallel Distributed Processing (Vol. 1)*, D. E. Rumelhart & J. L. McClelland (eds.). Cambridge, MA: MIT Press.
- Saito, K. & Nakano, R. (1988). Medical diagnostic expert system based on the PDP model. *Proceedings of the IEEE International Conference on Neural Networks* (pp. 255-262). IEEE Press.
- Scott, G., Shavlik, J., & Ray, W. (1992). Refining PID controllers using neural networks. *Neural Computation*, 4, 746-757.
- Sejnowski, T. J. & Rosenberg, C. (1987). Parallel networks that learn to pronounce English text. *Complex Systems*, 1, 145-168.
- Shavlik, J. W. & Towell, G. G. (1989). An approach to combining explanation-based and neural learning algorithms. *Connection Science*, 1, 233-255.

Combining Symbolic and Neural Learning

- Touretzky, D. S. (ed.) (1991). Special issue on connectionist approaches to language learning. *Machine Learning*, 7.
- Towell, G. G. (1992). *Symbolic Knowledge and Neural Networks: Insertion, Refinement, and Extraction*. Doctoral dissertation, Madison, WI: University of Wisconsin, Computer Sciences Department.
- Towell, G. G., Shavlik, J. W. & Noordewier, M. O. (1990). Refinement of approximately correct domain theories by knowledge-based neural networks. *Proceedings of the Eighth National Conference on Artificial Intelligence* (pp. 861-866), Boston: AAAI Press.
- Towell G. G. & Shavlik, J. W. (1992a). Using symbolic inductive learning to improve knowledge-based neural networks. *Proceedings of the Tenth National Conference on Artificial Intelligence* (pp. 177-182), San Jose, CA: AAAI Press.
- Towell, G. G. & Shavlik, J. W. (1992b). Interpretation of artificial neural networks: Mapping knowledge-based neural networks into rules. In *Advances in Neural Information Processing Systems (Vol. 4)*, J. Moody, S. Hanson, & R. Lippmann (eds.), San Mateo, CA: Morgan Kaufmann. (A longer version of this paper will appear in *Machine Learning* under the title "Extracting refined rules from knowledge-based neural networks.")
- Tresp, V., Hollatz, J., & Ahmad, S. (1993). Network structuring and training using rule-based knowledge. In *Advances in Neural Information Processing Systems (Vol. 5)*, S. Hanson, J. Cowans, & L. Giles (eds.), San Mateo, CA: Morgan Kaufmann.
- Uberbacher, E. C. & Mural, R. J. (1991). Locating protein coding regions in human DNA sequences by a multiple sensor - neural network approach. *Proceedings of the National Academy of Sciences*, 88, 11,261-11,265.
- Utgoff, P. E. (1988). Perceptron trees: A case study in hybrid concept representations. *Proceedings of the Seventh National Conference on Artificial Intelligence* (pp. 601-606). St. Paul, MN: Morgan Kaufmann.
- Weiss, S. & Kapouleas, I. (1989). An empirical comparison of pattern recognition, neural nets, and machine learning classification methods. *Proceedings of the Eleventh International Joint Conference on Artificial Intelligence* (pp. 781-786). Detroit: Morgan Kaufmann.

Appendix - An Application of KBANN

This appendix summarizes a study in which KBANN was used to produce a method for recognizing promoters in *E. coli* DNA (Towell et al., 1990); promoters are the sites where the process of "expressing" a gene to create a protein begins. Table 2 contains the initial domain theory used in the promoter recognition task (underscore, , indicates a position that can be filled by *any* DNA nucleotide). The first rule says that a promoter involves two subcategories: a `contact` and a `conformation` region. The second rule states that a `contact` involves two regions, while subsequent rules define alternative ways these regions can appear. The set of rules was derived in a straightforward fashion from the biological literature that describes the DNA sequences used for training (O'Neill, 1989).

The input features are 57 sequential DNA nucleotides. A special notation (the "@" symbol) is used to simplify specification of locations in the DNA sequence. The biological literature counts locations relative to the site where gene expression begins. Fifty nucleotides before and six following this location constitute an example. When a rule's antecedents refer to input features, they first state the starting location, then list a sequence that must follow. Hence, the second rule for `conformation` says that there must be an adenine (a) 45 nucleotides before the start site. Another adenine must be at position -44, then any two nucleotides can appear, and finally there must be a thymine (t) at location -41.

KBANN translates this domain theory into a neural network that has the topology shown in Figure 7. For clarity, the connections from `conformation` to the input units are not shown. Notice that in order to classify a subsequence as a promoter region, the intermediate concepts in

Table 2. An imperfect domain theory for DNA promoters.

```

promoter :- contact, conformation.      contact :- minus_35, minus_10.
minus_35 :- @-37 "cttgac".             minus_10 :- @-14 "tataat".
minus_35 :- @-36 "ttg_ca".             minus_10 :- @-13 "ta_a_t".
minus_35 :- @-36 "ttgaca".             minus_10 :- @-13 "tataat".
minus_35 :- @-36 "ttgac".              minus_10 :- @-12 "ta__t".
conformation :- @-47 "caa_tt_ac", @-22 "g__t_c", @-8 "gcgcc_cc".
conformation :- @-45 "aa__a"
conformation :- @-49 "a____t", @-27 "t____a__t_tg", @-1 "a".
conformation :- @-45 "a____a", @-28 "t__t_aa__t", @-4 "t".
    
```

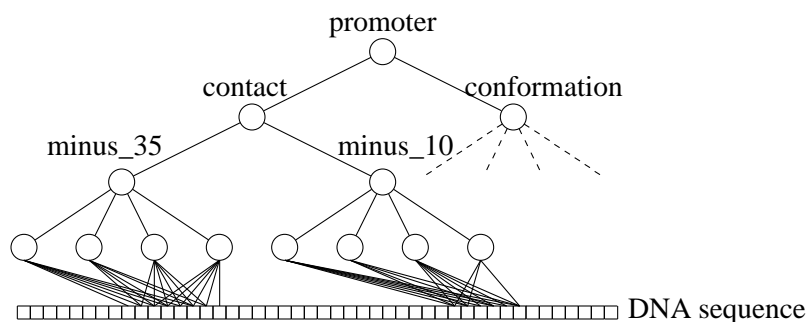


Figure 7. The initial neural network for promoter recognition that KBANN produces.

the domain theory must first be recognized. Also, recall that the algorithm adds additional lowly-weighted links (not shown) so that if additional sequence information is relevant, the algorithm can capture that information during backpropagation training.

Fifty-three sample promoters and 53 nonpromoter sequences were used to refine the initial neural network. The 53 sample promoters were obtained from a compilation produced by Hawley and McClure (1980). Negative training examples were derived by selecting contiguous substrings from a 1.5 kilobase sequence provided by Prof. T. Record of Wisconsin's Chemistry Department. By virtue of the fact that the fragment does not bind RNA polymerase (the protein that initiates gene expression), it is believed to not contain any promoter sites (Record, personal communication).

In order to get an estimate of how well the KBANN learned the concept of promoter, we used a standard experimental methodology called "leave-one-out" (also known as jackknife testing and cross-validation). This technique operates by training using $N-1$ examples, then

testing using the example left out. The procedure is repeated N times, so that each example is excluded once from the training set. The error rate is the number of errors on the single test cases, divided by N (here $N=106$).

Using the same methodology, we also applied two other learning algorithms: standard backpropagation and ID3 (Quinlan, 1986). In the standard backpropagation runs, we used the same number of hidden units (16) as in the network KBANN produced. Following convention, all of the input units were connected to each hidden unit and every hidden unit was connected to the output unit. All weights were randomly initialized to a number near zero. ID3 is a non-connectionist empirical learning algorithm. It uses training data to construct a decision tree for determining the category of an example. At each step, a new node is added to the decision tree by partitioning the training examples based on their value along a single, most-informative feature. The feature chosen is the one which, if chosen to partition the examples, maximizes a statistical measure of *information gain*.

Table 3 contains the average error rate on the training examples for the three learning algorithms (the original domain theory has an error rate of 50%, because it fails to exactly match any of the examples). In all cases, each algorithm correctly classified all members in the training sets. Hence, although each algorithm fully accounted for the training data, KBANN did a better job of *generalization*, in that its error rate on previously unseen examples was substantially lower. We also investigated a *nearest-neighbor* classification algorithm. The distance measure was the number of mismatched nucleotides. The k nearest neighbors were selected and the majority class of these neighbors chosen as the classification on the novel example. Testing values of k lead to a lowest error when $k=3$. Finally, the table reports the results of O'Neill's *ad hoc* approach (O'Neill, 1989), which is tailored to the specific task of promoter recognition.

Table 3. Error rates in the DNA promoter experiment.

KBANN	Standard Backpropagation	O'Neill	Nearest Neighbor	ID3
4/106	8/106	12/106	13/106	19/106

Towell (1992) reports the results of additional learning algorithms, including several that use symbolic machine learning techniques to refine domain theories; none of these alternates matched KBANN's generalization performance. See Towell and Shavlik (1992b) for the results of the NofM algorithm on the promoter problem. The promoter domain theory and training examples are available by anonymous ftp (`ics.uci.edu`) from the University of California at Irvine's Archive of Machine Learning Datasets and Domain Theories.