
Is Consistency Harmful?

William M. Spears
Naval Research Laboratory
Washington, D.C. 20375 USA
spears@aic.nrl.navy.mil

Diana F. Gordon
Naval Research Laboratory
Washington, D.C. 20375 USA
gordon@aic.nrl.navy.mil

Abstract

One of the major goals of most early concept learners was to find hypotheses that were perfectly consistent with the training data. It was believed that this goal would indirectly achieve a high degree of predictive accuracy on a set of test data. Later research has partially disproved this belief. However, the issue of consistency has not yet been resolved completely.

We examine the issue of consistency from a new perspective. To avoid overfitting the training data, a considerable number of current systems have sacrificed the goal of learning hypotheses that are perfectly consistent with the training instances by setting a new goal of hypothesis simplicity (Occam's razor). Instead of using simplicity as a goal, we have developed a novel approach that addresses consistency directly. In other words, our concept learner has the explicit goal of selecting the most appropriate degree of consistency with the training data.

We begin this paper by exploring concept learning with less than perfect consistency. Next, we describe a system that can adapt its degree of consistency in response to feedback about predictive accuracy on test data. Finally, we present the results of initial experiments that begin to address the question of how tightly hypotheses should fit the training data for different problems.

1 INTRODUCTION

Early studies in supervised concept learning made the implicit assumption that the best method for obtaining

high predictive accuracy on a test set is to find hypotheses that are perfectly consistent with respect to all examples in a training set (e.g., Michalski 1983). A positive hypothesis (i.e., a hypothesis intended to cover the positive examples) is 100% consistent with respect to a set of examples if it covers all positive examples and no negative examples in the set.

Perfect consistency was a goal for many years - until researchers began to examine more realistic databases that contained noisy, sparse data and unknown but possibly complex target concepts. To perform well on these databases, some systems sacrificed perfect consistency in favor of simplicity or other biases (Quinlan 1987; Michalski 1990). This achieved excellent results. Today, the issue of the ideal degree of consistency to use in a given situation (e.g., target concept and learning algorithm) is still unsettled. Some researchers, such as Angluin & Laird (1988) and Schaffer (1991), discuss the virtues of striving for 100% consistency. Shaffer (1991), for example, considers 100% consistency to be an appropriate bias for an "overwhelming majority" of situations. Other researchers, such as Quinlan (1987) and Michalski (1990), discuss the virtues of using a simplicity bias that sacrifices perfect consistency.

This paper examines the consistency issue from a new perspective. The simplicity bias is becoming prevalent in the concept learning literature. A simplicity bias typically satisfies two goals of the person who implements it: improved human understandability of the hypotheses, and improved predictive accuracy by avoiding overfitting the training data. For the sake of clarity and experimental precision, in this paper we adopt a novel approach that focuses *only* on the latter goal - we abandon the simplicity bias in favor of a bias that selects a consistency level. Each consistency level corresponds to a degree of fit to the training data, where 100% consistency implies the hypothesis fits the training data perfectly. Here, we present experiments that vary the consistency level, as

well as some initial answers to the question of when 100% consistency on the training data is best for achieving high predictive accuracy on the test data.

This paper also describes an adaptive approach to concept learning that views the consistency level as a bias that can be adjusted dynamically during learning. Because the goal of our learner is to improve its predictive accuracy, we decided to make that goal explicit by feeding predictive accuracy information back into our learner. Using this performance feedback, our learner selects the most appropriate consistency level to improve its predictive accuracy. This approach expresses a philosophy of a closed-loop feedback concept learner that has access to feedback about its ultimate goal. Surprisingly, such an approach is not frequently found in the literature. Exceptions include Breiman *et al.* (1984) and Michalski (1990).

In this paper, we examine one learning algorithm and consider the effects of varying the consistency level. Section 2 describes our concept learner, called the Genetic Algorithm Batch-Incremental Learner (GABIL), that we use in all experiments (De Jong *et al.* 1992). Section 3 describes a modified version of GABIL that can learn concepts with different levels of consistency. This section also presents experimental results that compare predictive accuracy with different consistency levels on both clean and noisy data and a variety of target concepts. Section 4 describes the closed-loop adaptive version of GABIL (AGABIL) that dynamically adjusts its consistency level to improve its predictions. Section 4 presents experimental results on the same data as Section 3, but this time using AGABIL. Section 5 relates this work to other research, and Section 6 states our conclusions and ideas for future research.

2 BACKGROUND: THE GABIL SYSTEM

GABIL is a supervised concept learning program based on the principles of Darwinian evolution and genetics (i.e., a genetic algorithm). In GABIL, Disjunctive Normal Form (DNF) hypotheses compete for survival, and reproduce according to their fitness with respect to a set of classified training instances (examples). Those hypotheses that are most fit survive and mate, producing new hypotheses via the application of genetic operators.

In GABIL, the concept of fitness is tied to that of consistency as follows:

$$fitness(hyp) = training_accuracy(hyp)$$

Accuracy refers to how well a hypothesis predicts the classification of a set of training examples. If a hypothesis predicts all the examples correctly, it is 100%

accurate. Similarly, if a hypothesis predicts one half of the examples correctly, it is 50% accurate. Training accuracy is equivalent to consistency. If a hypothesis is $N\%$ accurate, then it is $N\%$ consistent. For the sake of understandability, we write the fitness function:

$$fitness(hyp) = consistency(hyp)$$

to remind us that we reward those hypotheses that are more consistent. Note that simplicity plays no role in this fitness function.

A flowchart of GABIL is presented in Figure 1 (in all figures, ‘‘C’’ refers to consistency). GABIL is presented with two inputs: a set of training examples and a desired consistency level of 100% ($N = 100$ in Figure 1). GABIL returns as output a perfectly consistent hypothesis. This hypothesis is used to predict the classification of a new, previously unseen, example.

3 VARYING THE CONSISTENCY LEVEL

In this section, we examine the relationship between consistency level and predictive accuracy. To do this, we first modify GABIL so that the system can deliberately select hypotheses with less than 100% consistency. We then test the effects of varying the consistency level.

3.1 MODIFICATIONS TO GABIL

The initial version of GABIL always rewards those hypotheses that are most consistent. Suppose, however, that we wish to consider the effects of lower consistency

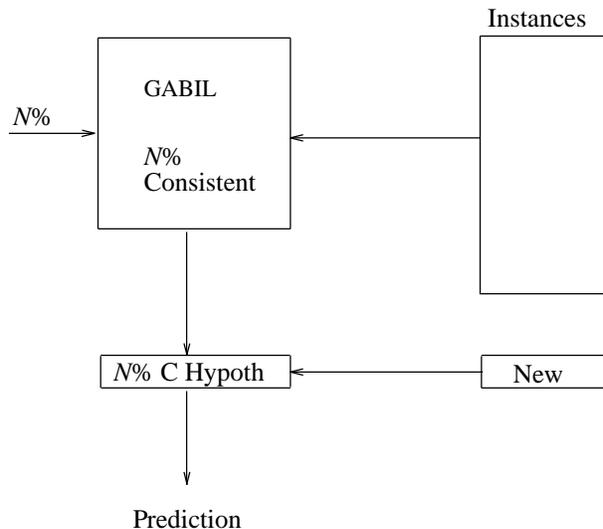


Figure 1: The flowchart for GABIL

levels. In other words, suppose we desire a consistency level of 90%. Then we wish to reward most highly those hypotheses that are closest to 90% in accuracy, and to reward less those hypotheses that are both more or less accurate. To do this, we modify the fitness function to:

$$fitness(hyp) = 1 - | consistency(hyp) - N |$$

in which N is the desired consistency level (see Figure 1). This function is maximized when the accuracy (consistency level) of the hypothesis matches N . The fitness is lower for those hypotheses that are both higher and lower than N in accuracy.

3.2 EXPERIMENTAL METHODOLOGY

We can now use GABIL to compare the effects of consistency level on predictive accuracy. Our experiments use a domain of artificial target concepts, which we call the $nDmC$ domain. In this domain, we have a four feature world, with four nominal values per feature (i.e., there are 256 instances in this domain). There are eight target concepts, that vary in complexity by increasing both the number of disjuncts and the number of relevant features (conjuncts) per disjunct. The number of disjuncts range from one to four, while the number of conjuncts is either one or three. Each target concept is labeled as $nDmC$, where n is the number of disjuncts and m is the number of conjuncts (see the Appendix for the definition of these target concepts). GABIL learns one target concept at a time.

Although GABIL is illustrated in Figure 1 as performing in batch mode, it is also capable of performing in a batch-incremental mode (i.e., batch mode is repeated for every new example). This allows us to generate incremental learning curves for the 256 instances in the $nDmC$ domain. Figures 3 - 6 depict a few representative learning curves. All learning curves are averaged over 10 independent runs for each target concept. For the sake of brevity, in our tables we present the global average of the predictive accuracy over each curve.

Since the issue of appropriate levels of consistency is intimately tied to that of noise, we examine both noise-free data and data with classification noise. This paper does not examine attribute noise, so we will refer to classification noise as simply ‘‘noise’’. We present results for noise-free data and data with 20% noise. We define $n\%$ noise such that each instance has a $n\%$ probability of receiving a random classification. Thus, 100% noise refers to the situation where the target concept is totally obscured. It is important to note that, because every instance is unique, an increase in noise is equivalent to increasing the target concept complexity,

and a perfectly consistent hypothesis is always possible. In this paper, we assume the source is noise. However, we also consider what our experimental results would imply if the source were instead increased target concept complexity.

3.3 RESULTS

As mentioned earlier, our motivation in this section is to examine the effects of consistency level on predictive accuracy. We ran GABIL with consistency levels of 100%, 90%, and 80% on the $nDmC$ target concepts. Tables 1 - 2 present the global averages of predictive accuracy for each target concept with 0% and 20% noise. We use $N\%$ to denote GABIL’s consistency level in the tables. A ‘‘*’’ highlights the winner (i.e., highest predictive accuracy) for each target concept. ‘‘ Δ ’’ is the difference in predictive accuracy between GABIL with 100% and 90% consistency. ‘‘Sig’’ is a two-tailed

Table 1: Effect of consistency level

TC	0% Noise				
	100%	90%	80%	Δ	Sig
1D1C	95.4*	85.3	76.5	+10.1	95%
1D3C	96.5*	88.1	79.0	+8.4	95%
2D1C	92.6*	81.1	71.7	+11.5	95%
2D3C	94.1*	88.5	81.5	+5.6	95%
3D1C	90.2*	77.0	69.2	+13.2	95%
3D3C	91.0*	87.8	79.5	+3.2	95%
4D1C	88.8*	75.2	65.8	+13.6	95%
4D3C	88.4*	86.6	79.3	+1.8	90%

Table 2: Effect of consistency level

TC	20% Noise				
	100%	90%	80%	Δ	Sig
1D1C	78.3	81.5*	74.8	-3.2	95%
1D3C	76.5	83.9*	78.7	-7.4	95%
2D1C	77.2	77.8*	69.9	-0.6	<80%
2D3C	75.8	82.8*	78.4	-7.0	95%
3D1C	77.6*	74.5	67.5	+3.1	90%
3D3C	75.1	80.3*	77.4	-5.2	95%
4D1C	77.2*	74.3	63.0	+2.9	95%
4D3C	73.8	78.9*	77.2	-5.1	90%

statistical test of significance for that difference.¹

Table 1 illustrates a virtue of 100% consistency. With 0% noise, GABIL with 100% consistency performs better than 90% and 80% consistency levels, on all target concepts. With 20% noise, however, the results are reversed. GABIL with 90% consistency outperforms 100% consistency on six of eight target concepts. Clearly, 100% consistency is a disadvantage in this situation.

In summary, it is certainly not the case that a particular consistency level is most appropriate for all target concepts and amounts of noise. Perfect consistency appears to be appropriate for some situations, and less than perfect consistency is appropriate for others. Therefore, it is natural to ask whether an adaptive mechanism can successfully determine an appropriate level of consistency. We address this issue in the following section.

4 ADAPTIVE CONSISTENCY LEVEL

Recall from Section 3 that GABIL can search for a hypothesis with a desired degree of consistency. In our previous section, this was manually controlled to examine the effect of consistency level on predictive accuracy. Suppose, however, that we could automatically determine an optimal level, *while* GABIL is learning a particular target concept. The advantages of such a mechanism are two-fold. First, we can analyze the adaptive mechanism to see what level of consistency it chooses for a particular target concept and level of noise. Second, this approach follows the valuable control theory philosophy for closed-loop feedback systems, i.e., if you wish to optimize predictive accuracy, this information should be available to the system. Practically speaking, the resulting system can be more robust because it can monitor its own performance.

4.1 MODIFICATIONS TO GABIL

We modified GABIL to create an adaptive system (AGABIL). AGABIL makes two passes over the training data before predicting the class of each new instance.² On the first pass, the training data is split into two sets, which we denote as A and B. Set A contains 3/4 of the data, while B contains the remaining 1/4. On the first pass, AGABIL searches for a perfectly consistent hypothesis with respect to set A. AGABIL also stores a small number of

¹ The significance test assumes nearly equal variances. We increased the rigor of the significance test whenever the sample variances differed according to the *F*-statistic.

² Our two-pass method is similar to the cross-validation method described in Breiman *et al.* (1984). CPU time prohibited the use of more than two passes in AGABIL.

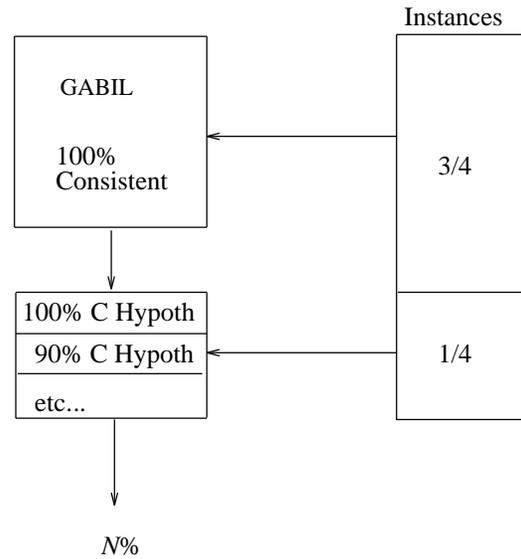


Figure 2: The first pass of AGABIL

less consistent hypotheses as it searches for a 100% consistent hypothesis. In our experiments, AGABIL stores one hypothesis in each of the four ranges: 70-79%, 80-89%, 90-99%, and 100%. AGABIL then compares the predictive accuracy of these stored hypotheses on set B. The consistency level of the best predicting hypothesis is chosen and recorded. The first pass is illustrated in Figure 2.

Next, a second pass over the complete training data is made to find a hypothesis that achieves the chosen consistency level over this data. The hypothesis resulting from the second pass (with the chosen consistency level) is then used to predict the class of a new (test) instance.³ Figure 1 illustrates the second pass. Note that in this adaptive system the consistency level is no longer set by the user, but is instead determined by a preliminary pass over the training data.

4.2 RESULTS

AGABIL was run on the *nDmC* domain, again with 0% and 20% noise. Tables 3 - 4 illustrate the results. In these tables, Δ is the difference in predictive accuracy between AGABIL and GABIL with 100% consistency (see Tables 1 - 2). Figures 3 - 6 show representative learning curves from which the predictive accuracy (denoted "PA" in our figures) averages are derived. In these figures, the solid

³ Experiments in which the resulting first pass hypothesis was used to predict the class of the new (test) instance were not as successful.

curve is the predictive performance of AGABIL, while the dotted curve is the predictive performance of 100% consistent GABIL.

When there is no noise, AGABIL performs well, nearly matching the performance of the best consistency level (100%) on the simpler concepts, and outperforming that consistency level on two of the more difficult concepts. When there is 20% noise, the results indicate quite strongly that the adaptive system can outperform GABIL with 100% consistency. In general, according to Tables 3 - 4, AGABIL performs slightly better in relation to 100% consistent GABIL as the number of conjuncts increases for a fixed number of disjuncts in the target concept. Tables 3 - 4 also show that when the adaptive system wins, the results tend to have a higher level of statistical significance than when it loses. Furthermore, AGABIL performs competitively with the best consistency level (see Tables 1 - 2).

Table 3: Performance of adaptive consistency level

TC	0% Noise		
	AGABIL	Δ	Sig
1D1C	94.4	-1.0	90%
1D3C	96.3	-0.2	<80%
2D1C	92.1	-0.5	<80%
2D3C	93.9	-0.2	<80%
3D1C	89.4	-0.8	<80%
3D3C	92.7	+1.7	90%
4D1C	88.1	-0.7	<80%
4D3C	90.5	+2.1	95%

Table 4: Performance of adaptive consistency level

TC	20% Noise		
	AGABIL	Δ	Sig
1D1C	81.9	+3.6	95%
1D3C	83.4	+6.9	95%
2D1C	79.5	+2.3	90%
2D3C	81.7	+5.9	95%
3D1C	76.8	-0.8	<80%
3D3C	81.0	+5.9	95%
4D1C	76.7	-0.5	<80%
4D3C	79.2	+5.4	95%

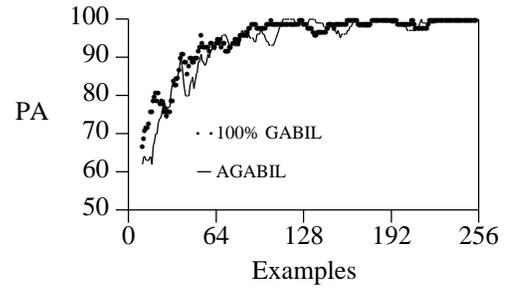


Figure 3: 1D1C - 0% noise

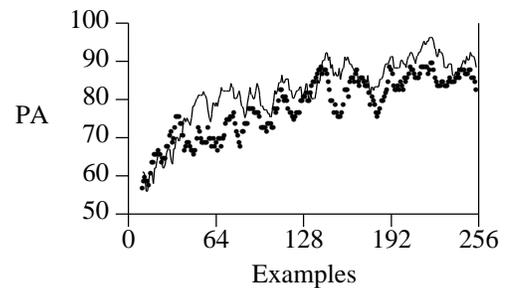


Figure 4: 1D1C - 20% noise

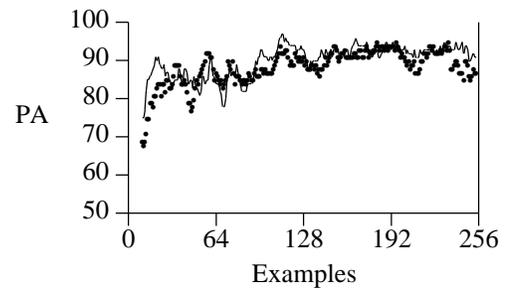


Figure 5: 4D3C - 0% noise

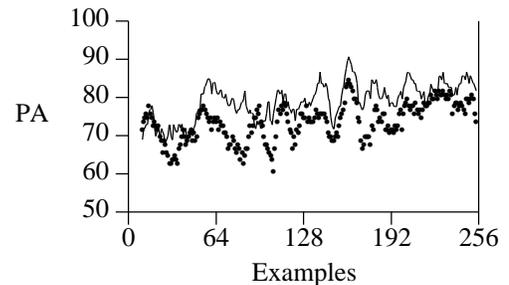


Figure 6: 4D3C - 20% noise

These results seem to indicate a trend. AGABIL shows better performance when the noise increases. Recall that we can equate an increase in noise with an increase in target concept complexity. Therefore, we conclude that the adaptive system performs better in relation to the 100% consistent system as noise or target complexity increases.

4.3 A POSSIBLE CAVEAT

The results in Section 4.2 indicate that lowering the consistency level below 100% is better when the target concept complexity increases. However, these results are particular to GABIL, which learns only the target concept, and not the negation of the target concept. One could argue that as the target concept becomes more complex, the negation of the target concept may become simpler. In these cases it might be reasonable to learn the negation of the target concept and aim for a consistency level of 100%. If this line of reasoning is correct, we would be entertaining the possibility that the important issue is *not* whether to strive for less than 100% consistency, but rather whether to learn the target concept or the negation of the target concept.

Since the complexity of an arbitrary target concept is not known beforehand, we need some measure to help us determine that complexity as the system runs. One possible measure is the ratio of positive to negative instances. As Schaffer (1991) indicates, the average target concept is simpler for those concepts with a preponderance of positive or negative examples. We examined this ratio for our *nDmC* domain and found that only 3D1C and 4D1C have more positive than negative examples. It is interesting to note that, according to the results in Section 4.2, 100% consistency produces better results than lower consistency levels on both 3D1C and 4D1C, regardless of the noise level. This result is true for these target concepts *only*. Therefore, the following heuristic appears to be a valid alternative to selecting the best consistency level:

IF $pos > neg$
 THEN learn TC with 100% consistency,
 ELSE learn $\neg TC$ with 100% consistency

where pos is the number of positive examples, neg is the number of negative examples, TC is the target concept, and $\neg TC$ is the negation of the target concept. An underlying assumption of this heuristic is that learning the target concept should produce better results than learning the negation of the target concept when there are more positive examples, and vice versa when there are more negative examples. Another underlying assumption is that 100% is the best consistency level for which a system should strive.

Although we have not had time to implement this heuristic within GABIL, we were able to compare the performance of GABIL with consistency levels of 100% and 90%, while learning the negation of the 3D1C and 4D1C target concepts.⁴

The results, which are shown in Table 5, are quite unexpected. First, learning the negation of the target concept is always better than learning the target concept, despite the preponderance of positive examples. Second, there is still evidence that a lower (90%) consistency level is more useful than 100% consistency as the noise increases and the target concept becomes more complex. These results clearly diminish the general usefulness of our heuristic. Therefore, we continue to stress the importance of adaptively adjusting the consistency level.

These results also suggest that it may be difficult to decide *a priori* which target concept for GABIL to learn. One possible solution is to let the system learn both the target concept and its negation simultaneously. However, there are a number of implementation issues that make this solution infeasible. A more intriguing solution is to let GABIL adaptively decide both its consistency level and the target concept it will learn, based on predictive performance. We will pursue this possibility in future implementations of GABIL.

Table 5: Learning the negation

	0% Noise		20% Noise	
	100%	90%	100%	90%
TC				
3D1C	90.2	77.0	77.6	74.5
\neg 3D1C	95.7*	85.0	77.6	80.8*
4D1C	88.8	75.2	77.2	74.3
\neg 4D1C	92.2*	84.0	78.4	81.6*

4.4 CHANGES IN CONSISTENCY LEVEL

One of the advantages of an adaptive mechanism is robustness. This advantage has been shown in Section 4.2. Another advantage is that the adaptive system can be monitored. Figures 7 - 10 illustrate how the consistency level (denoted " $N\%$ ") within AGABIL changes, for particular target concepts and levels of noise. These figures highlight some interesting points. First, the appropriate

⁴ The negation of 3D1C is a 1D3C concept, and the negation of 4D1C is a 1D4C concept, for this particular domain.

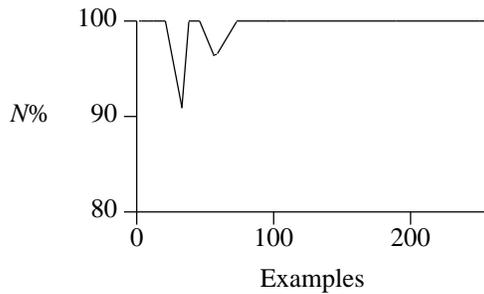


Figure 7: 1D1C - 0% noise

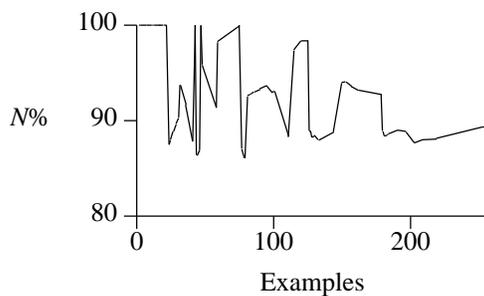


Figure 8: 1D1C - 20% noise

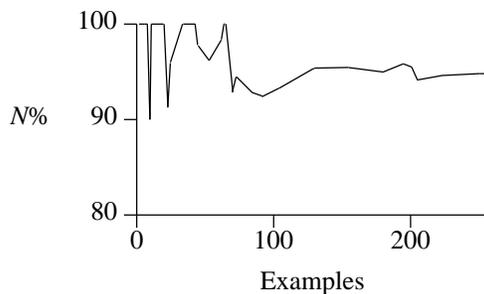


Figure 9: 4D3C - 0% noise

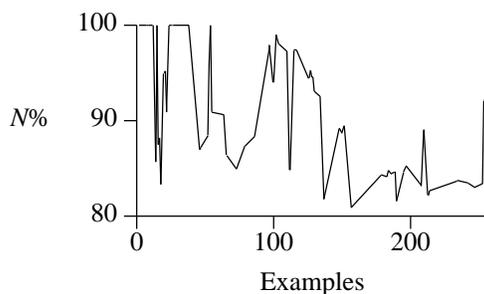


Figure 10: 4D3C - 20% noise

level of consistency is lower both for more complex target concepts and for greater levels of noise (which is analogous to a more complex target concept). Second,

the consistency level varies more when fewer examples are presented. This indicates that the adaptive mechanism is having some trouble early on, possibly due to insufficient sampling. Finally, and perhaps most interestingly, the consistency level usually drops as the number of examples increases. This behavior appears similar to that described by Fisher and Schlimmer (1988).

5 RELATED WORK

There have been many methods for handling noisy data, such as weighted hypotheses (Schlimmer & Granger 1986), Bayesian approaches (Buntine 1991), multiple version spaces (Mitchell 1978), and tree pruning (Quinlan 1987; Breiman *et al.* 1984). The goal of our research is to vary the consistency level to handle noisy data and complex concepts. No previous research has had precisely the same goal. The most closely related research investigates the effectiveness of a simplicity bias. This research is related because increased simplicity can result in a reduced consistency level.

Simplicity biases have been implemented with two of the most widely used hypothesis representations: decision trees and DNF hypotheses. Pruning (Quinlan 1987), which applies to decision trees, can reduce the consistency level because each decision tree branch that is pruned away may contain information to distinguish the classes of instances. After pruning, this information is lost. The removal of hypothesis disjuncts (Michalski 1990) is an effective method to increase the simplicity of DNF hypotheses. This method may sacrifice 100% consistency because the removed disjuncts may uniquely cover some of the training examples.

Breiman *et al.* (1984) and Quinlan (1987) have demonstrated that the simplicity bias is highly effective on a number of real-world domains, including domains containing noisy data. On the other hand, Schaffer (1991) claims that 100% consistency is usually better than selecting greater simplicity when trying to improve the predictive accuracy of a concept learner, even on data with classification noise. Schaffer draws this conclusion from a comparison of the CART system of Breiman *et al.* (1984), that has an adaptive strategy for selecting the best level to which to prune a decision tree, with a “naive” system that always maintains 100% consistency. The “naive” system keeps the full unpruned decision tree. Schaffer’s experiments indicate that CART will only outperform this “naive” system when the target concept is simple and there is little classification noise.

Some of our results, however, seem to conflict with those of Schaffer. For example, AGABIL usually performs better in relation to 100% consistent GABIL as target complexity increases. Also, AGABIL usually performs

better as classification noise increases. Further experimentation will be required to determine whether our results conflict because AGABIL's goal and CART's goal differ, or because other biases differ. The goal of our system is to find the appropriate consistency level, whereas the goal of the latter system is to find the appropriate simplicity level.⁵

Finally, if further experiments indicate that the difference in goals (i.e., selecting consistency versus selecting simplicity) accounts for the differences between our results and Schaffer's, then this might justify the advantage of using AGABIL's goal, rather than CART's, on real-world databases if human understandability is not an objective. The reason for preferring the goal of selecting consistency would be that AGABIL, which implements this goal, seems to perform better as the noise level or the target concept complexity increases.

6 CONCLUSIONS AND FUTURE WORK

In this paper, we have addressed the issue of finding an appropriate consistency level for improving predictive accuracy. Given a suite of target concepts that incrementally increase in complexity, and a corresponding set of training examples that vary in their level of noise, we identify the "best" consistency level for each case. The "best" consistency level is one that yields the highest predictive accuracy. We also describe a method for feeding the predictive accuracy information back into a learner to dynamically adjust the consistency level bias. Finally, we compare the performance of this adaptive system with a system that maintains 100% consistency over the training examples.

From these experiments, we have formed the following conclusions. First, lowering the consistency level seems to be more appropriate as the noise increases. Second, lowering the consistency level also seems to be more appropriate as target concept complexity increases. Finally, we have developed an adaptive concept learner that can select the best consistency level by using predictive accuracy feedback. This adaptive system is novel because it uses the predictive accuracy feedback to select a consistency level, rather than to select a simplicity level.

Future work will focus on three major directions. Our first direction will be to compare the lower consistency bias with the greater simplicity bias, in order to learn

when each bias is more effective. Another direction for future research will be to further test the generality of our results by rerunning our experiments using different systems and a wider variety of target concepts.

Our third direction for future research relates to the results in computational learning theory. Valiant (1984) has introduced the criterion of Probably Approximately Correct (PAC) identification of a target concept. Recently, a number of researchers have considered the computational feasibility of PAC identification in the context of noisy examples (e.g., Angluin & Laird 1988). However, they assume the strategy is to maximize consistency with the training sample. It would be interesting to also explore the computational feasibility of PAC identification assuming a strategy of lower consistency.

Acknowledgements

Special thanks to Ken De Jong, who made many major contributions to the GABIL project. We also thank the anonymous reviewers and the members of the Machine Learning Group at NCARAI for their useful comments on this research.

References

- Angluin, D. & Laird, P. (1988). Learning from Noisy Examples. *Machine Learning*, 2.
- Breiman, L., Friedman, J., Olshen, R., & Stone, C. (1984). *Classification and Regression Trees*. Belmont: Wadsworth.
- Buntine, W. (1991). Classifiers: A theoretical and empirical study. In *Proceedings of the Twelfth International Joint Conference on Artificial Intelligence*.
- De Jong, K., Spears, W., & Gordon, D. (1992). Using genetic algorithms for concept learning. To appear in *Machine Learning*.
- Fisher, D. & Schlimmer, D. (1988). Concept Simplification and Prediction Accuracy. In *Proceedings of the Fifth International Conference on Machine Learning*.
- Michalski, R. (1983). A theory and methodology of inductive learning. In R. Michalski, J. Carbonell, & T. Mitchell (Eds.), *Machine learning: An Artificial Intelligence Approach* (Vol. 1). Palo Alto: Tioga.
- Michalski, R. (1990). Learning flexible concepts: Fundamental ideas and a method based on two-tiered representation. In Y. Kodratoff, R. Michalski (Eds.), *Machine learning: An Artificial Intelligence Approach* (Vol. 3) San Mateo: Morgan

⁵ Schaffer (1991) also considers the effects that the hypothesis language bias and the ratio of positive to negative examples have on his conclusions. See Section 4.3 for our discussion of AGABIL and this ratio. We have not yet experimented with variations in the representational (hypothesis language) bias.

Kaufmann.

Mitchell, T. (1978). *Version spaces: An approach to concept learning*. Ph.D. thesis, Stanford University, Stanford, CA.

Quinlan, J. (1987). Simplifying decision trees. *International Journal of Man-Machine Studies*, 27.

Schaffer, C. (1991). Overfitting avoidance as bias. In *Proceedings of the Workshop on Evaluating and Changing Representation in Machine Learning* at IJCAI.

Schlimmer, J. & Granger, R. (1986). Incremental learning from noisy data. *Machine Learning*, 1.

Valiant, L. (1984). A theory of the learnable. *Communications of the ACM*, 27.

Appendix: $nDmC$ Target Concepts

This appendix fully describes the target concepts of the artificial domain. There are four features, denoted as F1, F2, F3, and F4. Each feature has four values $\{v1, v2, v3, v4\}$.

All the target concepts have the following form:

$$4DmC == d1 \vee d2 \vee d3 \vee d4$$

$$3DmC == d1 \vee d2 \vee d3$$

$$2DmC == d1 \vee d2$$

$$1DmC == d1$$

For the $nD3C$ target concepts we have:

$$d1 == (F1 = v1) \& (F2 = v1) \& (F3 = v1)$$

$$d2 == (F1 = v2) \& (F2 = v2) \& (F3 = v2)$$

$$d3 == (F1 = v3) \& (F2 = v3) \& (F3 = v3)$$

$$d4 == (F1 = v4) \& (F2 = v4) \& (F3 = v4)$$

Finally, we define the $nD1C$ target concepts:

$$d1 == (F1 = v1)$$

$$d2 == (F1 = v2)$$

$$d3 == (F1 = v3)$$

$$d4 == (F1 = v4)$$