

Flexibly Exploiting Prior Knowledge in Empirical Learning

Julio Ortega
Computer Science Department
Vanderbilt University
Nashville Tennessee 37235
U S A
juho@vu&e vanderbilt edu

Doug Fisher
Computer Science Department
Vanderbilt University
Nashville, Tennessee 37235
U S A
dfisher@vuse vanderbilt edu

Abstract

This paper presents a method to incorporate knowledge from possibly imperfect models and domain theories into inductive learning of decision trees for classification. The approach assumes that a model or domain theory reflects useful prior knowledge of the task. Thus the default bias should accept the model's predictions as accurate even in the face of somewhat contradictory data which may be unrepresentative or noisy. However our approach allows the system to abandon the model or domain theory, or portions thereof in the face of sufficiently contradictory data. In particular we use C4.5 to induce decision trees from data that have been augmented by 'model or domain-theory-derived features'. We weakly bias the system to select model-derived features during decision tree induction but this preference is not dogmatically applied. Our experiments vary imperfection in a model, the representativeness of data and the veracity with which model-derived features are preferred.

1 Introduction

When human expertise is nonexistent or very weak relative to a particular domain/task and when data is plentiful machine induction from data may be the only reasonable approach to task automation. In contrast, when expertise is strong, then encoding the expert's model or domain theory via traditional knowledge acquisition strategies may be the best approach. In fact, this human expertise may stem from induction over a much larger data sample than is available at the time task automation is undertaken.

In many cases, however, conditions are indeterminate as to whether sole reliance on machine induction or human expertise is most appropriate. Human expertise may not be 'perfect' and/or data may not be as plentiful as desired. In cases where some data is available and human expertise is less than perfect an advantageous strategy may be to exploit both in an appropriate way.

There is a growing body of work that combines model-based or domain-theory knowledge with empirical learning from data. Clark and Matwin [1993] assume that

an analyst-specified model mediates empirical learning - the rules derived from a machine-induction system are accepted as long as they do not contradict the biases found in the model. Evans and Fisher [1994] employ a similar strategy - a human analyst may specify weak rules (e.g. when printing-plant humidity is low, a certain kind of printing error known as banding is more likely, to occur). If inductively-derived rules indicate an opposite trend then the learning system's default strategy is to reject the rule derived through induction. In the case of both these approaches the model or domain theory is deemed correct in its characterization of the domain task though it may not be a very deep characterization. Inductive learning is used to flesh out rules that are consistent with the model (e.g. by selecting the particular numeric thresholds that distinguish high from moderate and low) or discovering rules relevant to a part of the domain space that are not addressed by the model or weak domain theory at all.

In the approaches above if the data contradicts the model then the implicit assumption is that the data are noisy or unrepresentative drawn from a very small subspace of the data. Other approaches known as *theory revision methods* [Ourston 1991] [Towell et al 1990] may give more credence to the data. In these systems contradictions result in a revision of the domain theory to bring it in line with the data. Drastal et al [1980] Rendell and Seshu [1990] and Ortega [1994] suggest an alternative strategy that loosely couples empirical learning and model-based reasoning. The data is augmented by features that are actually intermediate terms of the domain theory and which are deemed true of a datum by deductive application of the domain theory. Induction is then performed over this augmented data set. If domain-theory-derived features are included in rules derived inductively, then this suggests a rough consistency between the model and data. Model features may be viewed as somewhat better predictors than raw features because noise is mitigated. If model features are not referenced in a resultant classifier this may speak to imperfections in the model and/or this behavior may stem from an unrepresentative data sample. In both cases model-derived features may not look as informative as 'raw' features relative to the *available* data.

This paper describes a strategy that augments data with domain-theory derived features but unlike previ-

ous work we bias an adaptation of C4.5 [Quinlan 1993] to select domain-theory based features even when this conflicts somewhat with C4.5's original bias to select the most informative feature as computed over the data. The intent is to guard against the possibility of unrepresentative data. However, the domain-theory preference bias may be overridden if C4.5's original bias is sufficient opposed to the domain-theory preference bias. The intent here is to acknowledge that there may be some imperfections in the domain theory. Our experiments van imperfection in a model, the representativeness of data, and the veracity with which model-derived features are preferred.

2 Implementing a Flexible Domain-Theory Preference Bias

The approach described in this paper was motivated by our attempts to inductively build classifiers of faults of the Reaction Control System (RCS) of the Space Shuttle. A mixed qualitative/quantitative model for fault prediction was available [Robinson 1993] as well as simulated data representing system faults and normal behavior. For each available datum, the model was used to predict the fault. This prediction was added as a feature to the datum, as were various intermediate computations made by the model for the data point. The data points augmented in this way were then given to CA5, which constructed a classifier that predicted either a system's fault or normal operation. If the model were perfect, then we would expect that C4.5 would build a tree that only tested the model-based final prediction. Such a tree would indicate that if a new datum were encountered (represented by readings of various pressures and temperatures and other observables), then one should simply simulate the model on this datum and use the model-based final prediction. In the case of certain imperfections, a decision tree that tested various raw features, as well as various model-based features, might be constructed.

To our initial surprise, C4.5 consistently constructed trees that never or rarely referenced any model-based features. Rather than taking this as evidence of significant model imperfection or that the model added little or no information, above and beyond that implicit in the raw features, a NASA analyst familiar with this application indicated that the simulated data used for training was unrepresentative or skewed - it represented a very small subspace of the RCS description space.

This work motivated an approach that weakly biases our adaptation of C4.5 to select model-based features. In particular, for purposes of this paper, we assume a propositional domain theory used for classification that is acyclic and directed from the observable propositions to a final classification. A partial description of the perfect domain theory for the audiology domain used in our experiments is shown in Figure 1 as a tree. The domain theory is a set of rules, each one consisting of a set of conditions together with the classification predicted by the rule. In Figure 1, the antecedents of a rule are listed at the leaves of the tree. Each condition is an attribute-value pair (e.g., Air=pTofound). There may

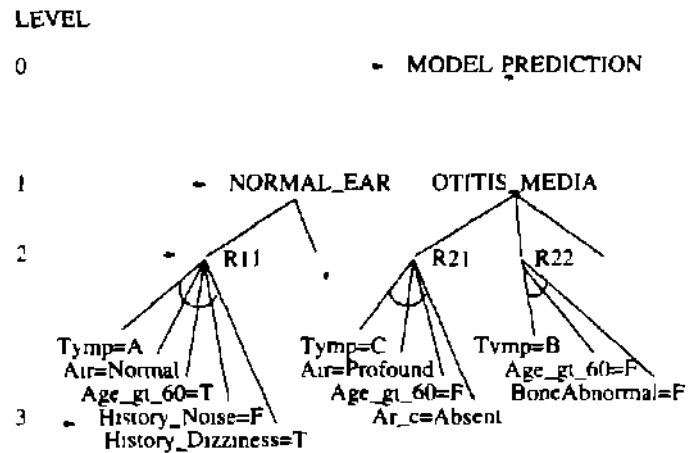


Figure 1 Levels in Audiology Theory

be several rules that predict a particular classification, as illustrated by the several possible rules leading to each classification (e.g., OTITIS MEDIA) in Figure 1.

We characterize model features extracted from a theory according to their distance in the theory hierarchy to the final model prediction.

- Level 0: Model prediction feature. We generate a single model prediction feature. Its value is the final prediction made by the rule interpreter (using the rules of the given theory).
- Level 1: Intermediate concept feature*. One of these features is generated for each possible classification in the theory. Each intermediate concept feature corresponds to a logical OR of the rules that predict particular classifications.
- Level 2: Rule features. Features of this type are generated for each rule in the theory. A rule feature follows from the logical AND of its antecedents.
- Level 3: Raw Features. Each rule's antecedent is a binary test as to whether an attribute takes a particular value on an example. Raw features, whether rule antecedents or not, are observable and are initially the exclusive means of representing data.

To bias C4.5 towards model features closer in the hierarchy to the final model prediction, we order features according to their level number from the model prediction feature through intermediate concept features to rule features, and raw features. At each step during induction, our variation of C4.5 chooses a feature of smallest level number unless a statistically-significant better feature (in terms of C4.5's information score) of larger level number is found. Hence, C4.5 will choose the model prediction feature unless sufficient evidence is present in the data to refute this choice.

Thus, we bias our inductive algorithm toward the model prediction feature and other features closer to it (of small level number). In a situation where we have a reasonably accurate model, and the available data is unrepresentative, we expect our model-biased method to work better than a default strategy of choosing the fea-

ture of highest information, value according to the available data (eg as in the standard C4.5). Nonetheless, if the data sufficiently contradicts the model the model bias can be abandoned and should we choose the model can be revised accordingly.

3 Using Domain-Theory Bias with Hypothesis Testing

The major difference between the original and our variation of C4.5 is the manner in which a feature is selected for each node, of a decision tree. C4.5 selects the feature with the highest information value according to the information gain ratio measure. Rather than selecting the feature with the highest information value outright our variation of C4.5 requires that this value be statistically significantly higher than the information value of all features preceding it in a feature preference ranking like that described in the previous section. Put in another way we select the highest feature in a preference ranking that has an information score not significantly worse than any feature lower in the preference ranking.

The above procedure is implemented by the function `SelectFeature(\ node)` shown in Figure 2 where F_p is the feature preference ranking, D is the set of training data associated with the current node, $info(f_j, D)$ is the value of C4.5's information measure for feature f_j when evaluated on the set of data D and F_j is the list of the features sorted in descending order according to this measure. `SelectFeature(\ node)` initially chooses the feature with highest information value (the first feature in F_j). However this feature is not accepted unless its information value is significantly higher than all features of higher preference according to the F_p ranking. If so the candidate feature is selected. Otherwise the highest preference feature found becomes the new candidate. The procedure is repeated until a significant difference is found or the F_j list is exhausted.

There is also a minor difference between the classification procedure of our system and the standard C4.5 algorithm for the situations where there is insufficient data to select a best for a particular node of the tree. As a purely data driven system the best C4.5 can do is to predict the most common class present in the current node. Instead since we assume our model is better than no information we use the prediction of our prior model.

The critical component of the function `SelectFeature` is the `SignificantlyBetter($f_{candidate}, f_{preference}, D$)` function shown in Figure 3. This function returns true if the information value of feature $f_{candidate}$ is estimated to be significantly higher than that of $f_{preference}$ according to a given level of statistical significance $sigLevel$. This is done by testing the null hypothesis that the difference between the information values of $f_{candidate}$ and $f_{preference}$ is zero. If this null hypothesis can be rejected with $1 - sigLevel$ confidence `SignificantlyBetter` concludes that $f_{candidate}$ is significantly better than $f_{preference}$.

¹In the current implementation the ranking is a total ordering. Features are sorted in ascending order according to level number. The ranking of features within a level number is arbitrary.

```

Given prior preference list  $F_p = f_1, f_2, \dots, f_n$ 

Function SelectFeature(\ node)

  Set  $D$  to set of observations in  $\ node$ 
  Create list of features  $F_j = f_{j1}, f_{j2}, \dots, f_{jn}$ 
    sorted in descending order according to value
    of  $info(f_j, D)$  eliminating any feature of null
    information value. In the case of nominal
    features this precludes the consideration of a
    feature used previously in the same path
  Set  $f_{candidate} = f_{j1}$ 
  While no significant difference has been found
  and there remain features to consider in  $F_j$ 
    Set  $f_{preference}$  to the first feature in  $F_j$  after
     $f_{candidate}$  that precedes  $f_{candidate}$  in  $F_j$ 
    Eliminate all other features
    between  $f_{candidate}$  and  $f_{preference}$  in  $F_j$ 
    from consideration
    If Not(SignificantlyBetter( $f_{candidate}, f_{preference}, D$ ))
      Set  $f_{candidate} = f_{preference}$ 
    EndIf
  EndWhile
  Return( $f_{candidate}$ )
  
```

Figure 2. Function `SelectFeature`

If the form of the probability distribution associated with C4.5's information measure is known and its parameters can be calculated then traditional statistical theory can be used to test significance. This could be done for the information gain measure since Musick et al. [Musick et al. 1993] prove that this measure is normally distributed and provide explicit formulas for the parameters of this distribution. However the form of the distribution for the default measure used in C4.5, information gain ratio, is not known. Fortunately *Bootstrap Methods*, [Lifson and Giong 1983] allow for estimates of significance levels of arbitrary statistics when the form and parameters of the underlying distribution are not known [Moreen 1989]. This is the method implemented in the function `SignificantlyBetter`.

In Lifson's Bootstrap methods an unknown complete population P is estimated by repeated uniform subsampling with replacement from an available sample D of P . From D we obtain a set of bootstrap subsamples $P_B = \{D_1, D_2, \dots, D_{N_B}\}$ where N_B is a prescribed number of subsamples. Each D_i (with $1 \leq i \leq N_B$) is very likely to contain some duplicate- and be missing some observations from D with the result that the values of $info(f_j, D_i)$ for each feature f_j will likely be different on each bootstrap subsample D_i . Under some additional assumptions we then proceed as if the bootstrap samples were obtained from the actual population P .

`Significantly Better` uses two different bootstrap methods described by Moreen [1989]: the Normal Approximation Method, and the Shift Method. Figure 3 shows the computation of some quantities used in the above methods. $Diff$ is the difference in information value between $f_{candidate}$ and $f_{preference}$ computed on the set of

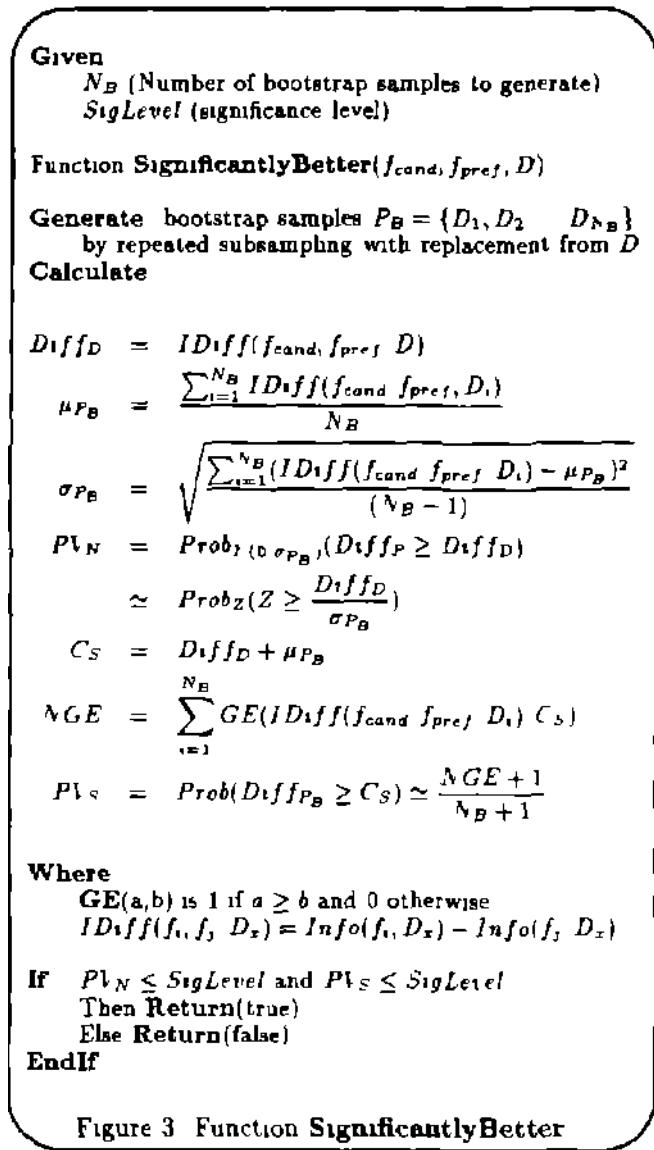


Figure 3 Function SignificantlyBetter

training data D up p_B , the mean of our statistic of interest (difference in the information value between f_{cand} and f_{pref}) over the the bootstrap samples and σ_{P_B} the standard deviation of our statistic of interest over the bootstrap samples

The Normal Approximation Method operates under the assumption that the sampling distribution of the statistic of interest (difference in the information value of f_{cand} and f_{pref}) under the null hypothesis (no difference) is normally distributed with mean zero and same variance as in the bootstrap samples. This assumption is used to calculate PV_N , the probability under the Normal/assumption that a value of our statistic higher than or equal to $Diff_D$ could have been obtained by chance. To calculate PV_N , we use σ_{P_B} and the probability function/tables for a standard normal distribution.

The Shift Method assumes that the sampling distribution of the statistic of interest on the complete population P has the same shape, but different mean, than the sampling distribution over the bootstrap samples P_B . To

determine the corresponding PVs , the probability under the *Shift* assumption that a value equal or higher than $Diff_D$ could have been obtained by chance we count the number of times that the value of the statistic on bootstrap samples is higher than a shift criterion ($C_S = Diff_D + \mu_{P_B}$), and divide that count by the number of subsamples N_B .

We only decide that the feature f_{cand} is significantly better than f_{pref} if it is significantly better according to both the Normal Approximation Method and the Shift Method. SignificantlyBetter is computationally quite expensive. However, during the selection of most features this needs to be done very few times. If the feature with the highest initial information value is the feature of highest preference, SignificantlyBetter never needs to be computed. When other features are initially selected, only the features with higher preference are checked. As soon as one significant difference is computed, no other significance computation is necessary. For the sake of efficiency, the case where two or more insignificant differences could account for a single significant difference is not considered. Our interest is not in precise computations of significance, but rather the qualitative effect of significance testing on the selection of attributes in CA 5 while retaining a reasonable level of efficiency.

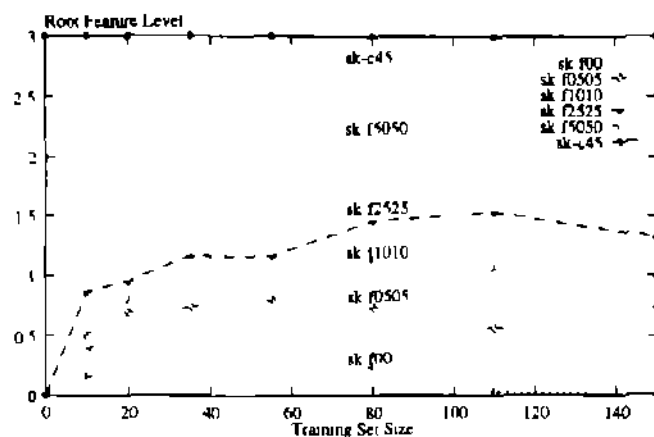
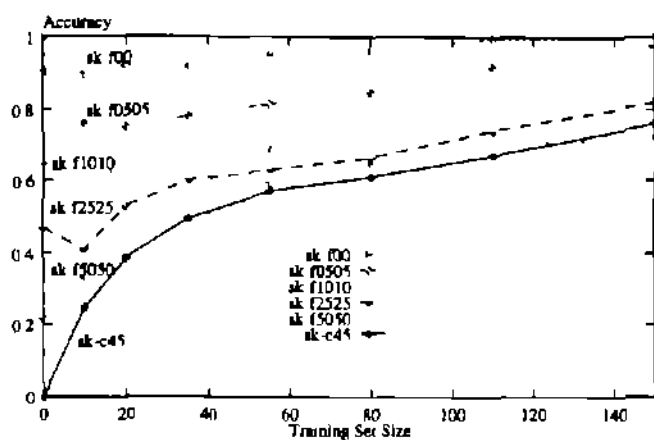
4 Experimental Design

To test our approach, we conducted experiments with the audiology dataset from the UCI (University of California at Irvine) Machine Learning repository. This database contains 226 examples from 24 classes. Each example is described by 67 discrete features.

For each of the 30 learning trials of our experiments, a test set of 76 examples and a disjoint training set of 150 examples were randomly and uniformly selected. Because we want to test the robustness of various strategies in the face of unrepresentative data, we sorted the training examples according to their Euclidean distance from a randomly selected datum. The training set was further divided into subsets which contain the first 10, 20, 35, 55, 80, 110, and 150 (i.e., all) examples of the sorted training set. Decision trees are learned for each of these subsets of training data.

The dependent variables of interest are predictive accuracy and the level (as illustrated in Figure 1) of the root feature of the decision trees learned under different conditions. Due to the recursive nature of decision tree induction, we expect that the tendencies observed at the root can be extrapolated to other nodes of the tree.

The independent variables are the size of the training set, the significance level used for hypothesis testing in our variation of CA 5, and the degree of model imperfection. Note that while varying the size of the training set, we are also varying its degree of skewedness, because training data are ordered based on Euclidean distance, small samples tend to be drawn from a small portion of the data description space, the larger the training dataset, the higher the proportion of the complete data set it covers, and thus the more representative it becomes. For the largest data set of 150 examples, all skewedness disappears since all 150 examples were ran-



domly chosen from the complete data set. We present results from skewed sampling as this tends to represent worst case conditions for our learning system. We have also experimented with traditional sampling (for all training sizes) thus allowing us to better tease apart the influence of skew and training set size though this paper does not elaborate on this issue. Significance levels of 50%, 25%, 10% and 1% are varied to indicate increasing confidence in the quality of our models.

We follow Mooney's approach [Mooney 1993] for generating theories of varying degrees of imperfection. A perfect theory, i.e., one that correctly classifies 100% of all audiology examples was first constructed by running C4.5 on the complete data set of audiology examples with all pruning disabled. This theory (named *f00*) contained 86 rules with an average of 7.79 antecedents per rule. Imperfect theories (named *f55*, *f1010*, *f2525* and *f5050*) were generated by randomly adding and then randomly deleting a percent of all conditions from the rules of the perfect theory (a corresponding 5%, 10%, 25% and 50%). This results in contaminated theories with errors both of omission and of commission. The accuracies over the complete data set of the imperfect theories *f55*, *f1010*, *f2525* and *f5050* are 89%, 65%, 46% and 21%, respectively. For comparison, the accuracy obtained if

we always predict the most frequent class in the dataset is 21%.

5 Experimental Results

Figure 4.2 shows results from a baseline study. The curve labeled *sk-c45* shows the results of standard C4.5 on skewed trials with 'raw' features only. The rest of the curves in this figure show the accuracy and root feature level for C4.5 when model features are added to the description of the data, but no preference ordering or hypothesis testing is done. We can see that accuracy improves significantly (over *sk-c45*) when a domain theory is exploited even for a low quality theory such as *f5050*. These results compare favorably to other systems tested on this domain [Mooney 1993] [Ourston 1991].

An interesting fact illustrated in Figure 4 is that the number of training examples required for C4.5 with a domain theory (*sk-fxx* curves) to produce an accuracy equivalent to the corresponding theory alone seems to be inversely proportional to the quality of the theory, the two extremes being *f5050* (0 training examples required to reach 21% accuracy) and *f00* (all training examples required to reach 100% accuracy). Ideally however (the accuracy of our system should be equal or better than the accuracy of the model alone, or the C4.5) learning algorithm alone. This only occurs for large enough or representative enough training data sets. This behavior is not unique to our system; it can be observed in the published learning curves shown for some systems that combine analytical and empirical learning [Pazzani and Kibler 1992] [Ourston 1991]. As we will see significance testing of ranked features appears to mitigate this undesirable behavior.

Figure 4 gives a good indication of the type of features selected with theories of varying quality. Standard C4.5 can only access raw features (level 3). C4.5 with features from a perfect model (i.e. *f00*) chooses almost exclusively (but not always) the model prediction feature (level 0). With lower quality theories C4.5 gradually chooses features of greater level numbers. However for the perfect model, there should be no reason to choose any feature other than the model prediction feature. This does not happen due in part to a known bias of the information gain ratio against features with many values (In the audiology domain the model prediction feature has 24 possible values, other model features are binary and raw features have few values). However as we will see next this problematic bias can be mitigated with the use of significance testing.

Figure 5 shows the effects of significance testing of ranked features when C4.5 is augmented with features from a perfect model. Rather than the 150 examples

2 Tables corresponding to the graphs in this paper on learning means and standard deviations can be found at <http://www.vuse.vanderbilt.edu/~dfisher/iech/reports/ijcai95-tables.pa>

3 Here we use standard C4.5 for learning. For testing we use the classification procedure described in section 2.1, i.e. predicting with the model rather than the most frequent class at leaves of the decision tree where there is insufficient data for further decomposition.

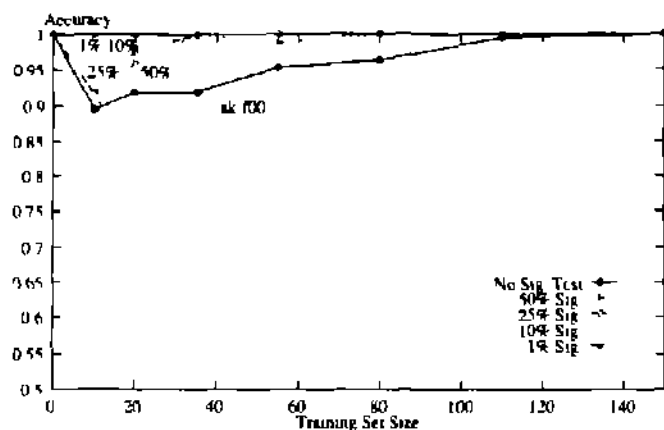


Figure 5 Average accuracy and root feature level of decision trees learned with features from a perfect model under varying levels of significance

needed before hypothesis testing of ranked features results in 100% accuracy with only 80 examples when using a 50% significance level or with just 85 examples when using more strict significance levels (25% 10% or 1%) This figure also shows how the average root feature level is gradually reduced with stricter significance levels

Perfect domain theories are an interesting boundary case, but most interesting theories are imperfect The graphs of Figure b illustrate how stricter levels of significance achieve our objective of biasing C45 toward features of smaller level number using the *f1010* domain theory This figure also shows that the accuracy obtained with the *f1010* imperfect theory improves consistently with stricter levels of significance testing (50% 40% 10%) for any size of the training set, including the complete training set of 150 examples In addition, while with no significance testing (or the almost equivalent significance testing at the 50% level) at least 55 examples are needed to obtain better accuracies than the *f1010* theory alone with stricter levels of significance testing only 10 examples are required For the 1% significance level, accuracy, is better than other significance levels when the training sets contains less than 110 examples, and worse when the training set contains more

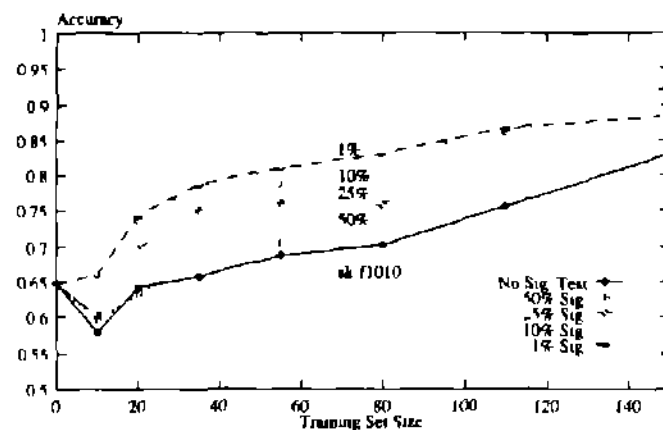


Figure b Average accuracy and root feature level of decision trees learned with features from the *f1010* imperfect model under varying levels of significance

than 110 examples Thus at least in this domain there is a breakeven point for significance levels that depends both on the quality of the model and the size of the training set after which stricter significance levels are detrimental In our experiments the size of the training sets corresponds to increasing quality in the available data in the sense that they are better representatives of the complete population both due to the sheer amount of data and due to the fact that the skewedness effect we introduced in the training set tends to diminish as the size of the training sets are increased

For lower quality theories similar behavior is observed with increasingly strict significance levels, accuracy improves when the training set contains few examples, and decreases with training sets that contain many examples Thus, for every combination of model quality (i.e. amount of contamination in the model) and data quality (level of skewedness and number of training examples) there is an optimal level of significance between the extremes of 50% and 0% in this domain However the choice of a beneficial but perhaps non-optimal significance level is not difficult Significance levels only seem to become detrimental for large data sets when we use significance levels stricter than 10% Values between

50% and 10% always seem to improve accuracy (perhaps non-optimally) for any size of the training set. Thus, expert intuition about the trustworthiness of an existing model with respect to the available data can be incorporated into our learning algorithm to obtain additional performance improvements.

6 Concluding Remarks

In this paper we address a situation that we believe to lie of practical interest: learning whenever we have a model believed to be of good quality but imperfect nevertheless, together with a set of data of unknown representativeness and quality. We present a method that attempts to take advantage of both the model and the data, plus our prior knowledge about the quality of the model. Our method biases empirical learning in a flexible manner such that model-based features, or more generally preferred features based on some a priori determined preference ordering, are selected unless sufficient refuting evidence appears in the data. The amount of evidence required is determined by statistical significance and is set by the user according to his/her confidence in the quality of the available model.

Our experimental results show that when features generated from a model are simply added to the description of the data, accuracy is increased to a degree proportional to the quality of the model. However, some problems with this simple approach are illustrated by the fact that perfect models only result in perfect accuracy with large or very representative sets of training examples; if significance testing with a preference ordering is used with a perfect model, our system becomes more robust in the presence of skewed data; few examples are then needed to obtain perfect accuracy. Further, with imperfect models of good quality, we obtain additional increases in accuracy for any number of training examples.

Although significance testing has been used previously in machine learning methods, such as for the pruning of decision trees [Quinlan 1986], our use of this concept for flexibly introducing prior knowledge bias in empirical learning seems to be novel.

Acknowledgements

This research was supported by a grant from NASA Ames Research Center (NAG 2-834) to Doug Fisher. We thank Deepak Kulkarni and Peter Robinson for early discussion on initial strategies that led to those described here.

References

- [Clark and Malum 1993] P. Clark and S. Matwin. Using qualitative models to guide inductive learning. In Proceedings of the Tenth International Conference on Machine Learning, pages 49-5b. Amherst, MA, 1993.
- [Drashtal et al., 1989] G. Drashtal, G. Czako, and S. Raatz. Induction in an abstraction space: a form of constructive induction. In Proceedings of the Eleventh International Joint Conference on Artificial Intelligence, pages 708-712. Detroit, MI, 1989.
- [Efron and Gong, 1983] Bradley Efron and Gail Gong. A leisurely look at the bootstrap, the jackknife, and cross-validation. *The American Statistician* 37(1): 36-48, 1983.
- [Evans and Fisher, 1994] R. R. Evans and D. Fisher. Overcoming process delays with decision tree induction. *IEEE Expert*, 9(1): 60-66, 1994.
- [Mooney 1993] R. J. Mooney. Induction over the unexplained: Using overly-general theories to aid concept learning. *Machine Learning* 10: 79-110, 1993.
- [Musick et al., 1993] Ron Musick, Jason (allell) and Stuart Russell. Decision theoretic subsampling for induction on large databases. In Proceedings of the Tenth International Conference on Machine Learning, pages 212-219. Amherst, MA, 1993.
- [Noreen 1989] Eric W. Noreen. *Computer Intensive Methods for Testing Hypotheses: An Introduction*. John Wiley & Sons, New York, NY, 1989.
- [Ortega 1994] J. Ortega. Making the most of what you've got: using models and data to improve learning rate and prediction accuracy. Technical Report (S-94-01). Computer Science Dept., Vanderbilt University. Abstract appears in Proceedings of the Twelfth National Conference on Artificial Intelligence, p. 1483. Seattle, WA, 1994.
- [Ortega In Preparation] J. Ortega. Making the Most of What YOU Got using Models and Data to Improve Prediction Accuracy. PhD thesis, Vanderbilt University, Nashville, TN. In Preparation.
- [Ourston 1991] D. Ourston. *Inductive Explanation-Based and Empirical Methods in Theory Revision*. PhD thesis, University of Texas, Austin, TX, 1991.
- [Pazzani and Kibler 1992] M. Pazzani and D. Kibler. The utility of knowledge in inductive learning. *Machine Learning* 9(1): 37-94, 1992.
- [Quinlan 1986] J. R. Quinlan. Induction of decision trees. *Machine Learning* 1: 81-106, 1986.
- [Quinlan 1991] J. R. Quinlan. *4.5 Programs for Machine Learning*. Morgan Kaufmann, San Mateo, CA, 1993.
- [Rendell and Seshu 1990] Larry Rendell and Raj Seshu. Learning hard concepts through constructive induction: framework and rationale. *Computational Intelligence* 6(5): 247-270, 1990.
- [Robinson 1993] Peter Robinson. Automated fault diagnosis algorithms for the reaction control system of the space shuttle. Technical Report FIA-93-05. NASA Ames AI Research Center, 1993.
- [Towell et al. 1990] G. G. Towell, J. Shavlik, and M. O. Soorwedier. Refinement of approximate domain theories by knowledge-based neural networks. In Proceedings of the Eighth National Conference on Artificial Intelligence, pages 861-86b. Boston, MA, 1990.