

# AqBC: A Multistrategy Approach for Constructive Induction

Seok Won Lee and Janusz Wnek\*

Machine Learning and Inference Laboratory, M.S. 4A5

George Mason University

4400 University Drive

Fairfax, VA 22030-4444

{swlee, jwnek}@aic.gmu.edu

## ABSTRACT

In order to obtain potentially interesting patterns and relations from large, distributed, heterogeneous databases, it is essential to employ an intelligent and automated *KDD (Knowledge Discovery in Databases) process*. One of the most important methodologies is an integration of diverse learning strategies that cooperatively performs a variety of techniques and achieves high quality knowledge. *AqBC* is a multistrategy knowledge discovery approach that combines *supervised inductive learning* and *unsupervised Bayesian classification*. This study investigates creating a more suitable knowledge representation space with the aid of unsupervised Bayesian classification system, *AutoClass*. *AutoClass* discovers interesting patterns from databases. Via constructive induction, these patterns modify the knowledge representation space so that the robust inductive learning system, *AQ15c*, learns useful concept descriptions of a taxonomy. *AqBC* applied to two different sample problems yields not only simple but also meaningful knowledge due to the systems that implement its parent approaches. *AqBC*'s good performance appears to be due to its integration of reliable unsupervised Bayesian classification, constructive induction and rule induction, and not to the presence of any component alone.

## 1. INTRODUCTION

The explosive growth of large scale databases far exceeds our ability to analyze them, requiring a new approach for intelligent and automated knowledge discovery [4]. We present *AqBC* [7], a multistrategy knowledge discovery approach to concept learning. *AqBC* extracts new knowledge, determines meaningful descriptions and applies the newly acquired knowledge in supervised learning. These descriptions and knowledge grow out of patterns identified by *AqBC*. A clustering method using unsupervised Bayesian classification, generates the newly organized knowledge, while a supervised inductive rule learning system generalizes the descriptions and expresses them in variable valued logic. These new concepts expand the knowledge representation space for the supervised inductive learning system.

The system employs constructive induction to create and enhance the knowledge representation space with the aid of the unsupervised Bayesian classifier, *AutoClass* [2]. *AutoClass* provides a maximum posterior probability grouping objects into classes. The constructed classes define abstract concepts, with descriptions learned from class members using the inductive learning system, *AQ15c* [13]. The abstract concept descriptions are then used to improve and expand the original representation space. This expanded representation space serves as a final setting for supervised concept learning of any attribute from the original examples by employing *AQ15c*. The multiple stage concept learning has the following properties:

- The task of inferring a set of classes and class descriptions that best fit and explain a given data set is placed on a firm theoretical foundation using Bayesian statistics.
- The abstract concept descriptions learned in the first stage can illustrate and associate the corresponding concept descriptions learned in the second stage, which generates a set of simple descriptive rules. This way, the hierarchical hypotheses structures discovered from the nested classifications provide valuable information that cannot be obtained from either system alone.

The first MONK problem, MONK1, [12] and US census data have been used for experimentation. The diagrammatic visualization system, *DIAV* [15] graphically interprets the knowledge representation spaces and shows the changes in the representation space caused by constructive induction. In this paper, we show that the newly created knowledge facilitates classification and, in turn, problem solving that employs classification or pattern recognition in large databases.

## 2. THE SIGNIFICANCE OF IMPROVING THE KNOWLEDGE REPRESENTATION SPACE

*Constructive Induction (CI)* is a concept proposed in the field of inductive concept learning [8] to solve learning problems in which the original representation space is inadequate for the problem at hand and needs to be improved in order to correctly formulate the knowledge to be learned. In other words, constructive induction hypothesizes new knowledge using a

---

\* Also with Science Applications International Corp., Tysons Corner, VA.

search process. In our study, we search for the best representation space transformation by applying the unsupervised Bayesian classifier, AutoClass.

CI is based on the idea that the quality of the knowledge representation space is the most significant factor in concept learning. If the representation space is of high quality, i.e. the chosen attributes/descriptive terms are highly relevant to the problem at hand, learning will be relatively easy and will likely produce easily understood hypotheses with high predictive accuracy. If the quality of the representation space is low, (i.e. the attributes are less relevant to the problem) learning will be complex and no method may be able to produce good hypotheses. CI searches for patterns in data, learned hypotheses, and knowledge from experts, using them to create a new knowledge representation space [14].

In order to find an appropriate representation space, clustering is an important way of summarizing and explaining data. A clustering system accepts a set of object descriptions (events, observations, facts) and produces a classification scheme over the observations. This system does not require an “oracle” to preclassify objects, but instead uses an evaluation function to discover classes that provide “good” conceptual descriptions. In our AqBC approach, the unsupervised Bayesian classifier, AutoClass, plays the clustering role and searches for the best model or classifications. Its evaluation function is based on Bayesian statistics. Once AutoClass finds the best model, AqBC creates a new attribute called *class* and augments the original data set with it. This new set of attributes is then passed to the second phase, which employs the supervised inductive learning system, AQ15c. The supervised learning should benefit from this new knowledge representation space for learning target concepts and produce more accurate descriptions of the target concepts.

To prove the quality of effectiveness of this approach and illustrate its properties, we applied AqBC to MONK1. The original representation space with the training examples denoted by “+” and “-” is shown in Figure 2A using DIAV.

Like most real world problems, the initial representation space is very disordered and ill-structured. In the improved representation space (Figure 2B), the examples representing the target concept are better ordered and therefore easier to generalize. The positive examples are properly grouped in the middle value of the constructed attribute, while the negative ones group together in the first and last values. The descriptive rule sets learned by the AqBC approach are shown in Figure 1.

In this case, **x7** is a new attribute that represents the clusters, *class0*, *class1* and *class2*, created by AutoClass. This new attribute augments the original attribute set and significantly changes the shape of the representation space. The resulting space can be divided by a single three-valued attribute, so the rules describing the concepts are trivial.

The combined approach successfully captures these meaningful subconcepts and thus successfully solves the problem (100% accuracy on both the training and test sets). Since most real world problems do not provide an “oracle” to guide the learning task, an unsupervised classifier is an attractive device to modify the representation space. We will

discuss in later sections how AqBC addresses this issue in an application to US census data.

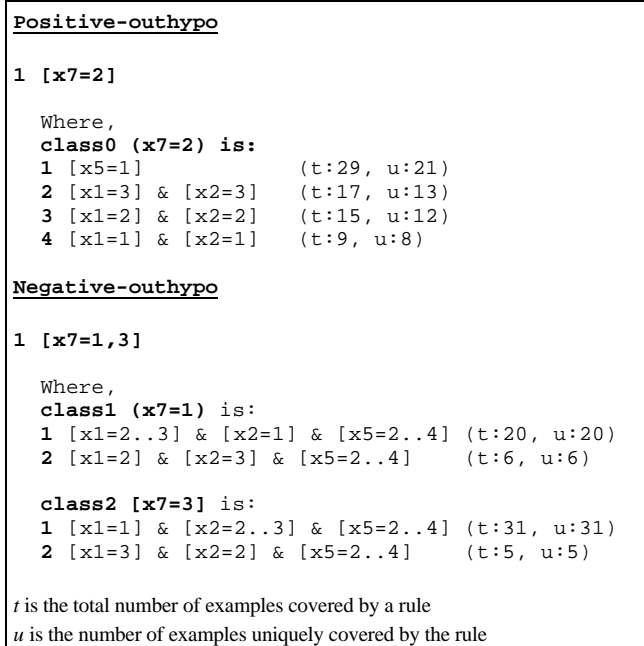


Figure 1. Rules learned by AqBC approach (MONK1 problem)

### 3. AQBC: A MULTISTRATEGY APPROACH FOR CONSTRUCTIVE INDUCTION-BASED KNOWLEDGE DISCOVERY

#### AQ15c Inductive Learning System

AQ15c [13] is a C-language reimplement of AQ15 [11] with significant performance improvement and dynamic allocation of representation spaces. AQ-family of inductive learning programs implements the STAR method of inductive learning [9].

AQ15c learns decision rules for a given set of decision classes from examples. When learning rules, AQ15c uses 1) background knowledge in the form of rules (input hypotheses), 2) the definition of descriptors and their types and 3) a rule preference criterion that evaluates competing candidate hypotheses. Each training example characterizes an object, and its class label specifies the correct decision associated with that object. The generated decision rules are expressed as symbolic descriptions involving relations between objects’ attribute values. The program performs a heuristic search through a space of logical expressions, until it finds a decision rule that best satisfies the preference criterion while covering all positive examples, but no negative examples.

Each decision rule is described by one or more conditions, all of which must be met for the rule to apply. A *condition* is a relational statement. A *rule* is a conjunction of conditions. A *hypothesis* is a disjunction of rules that together describe a concept.

The following is an example of a hypothesis consisting of two rules:

**Flag-outhypo**

1 [color = red, white, blue] & [stripes = 13] & [stars = 50]

2 [color = red, white, blue] & [stripes = 3]

A hypothesis is satisfied if any of its rules are satisfied, while a rule is satisfied if all of its conditions are satisfied. A condition is satisfied if the term takes one of the values in the reference. The hypothesis shown in the above example can be interpreted as follows: An object is a flag if: its color is red, white, or blue, and it has 13 stripes and 50 stars, **OR** its color is red, white, or blue, and it has 3 stripes.

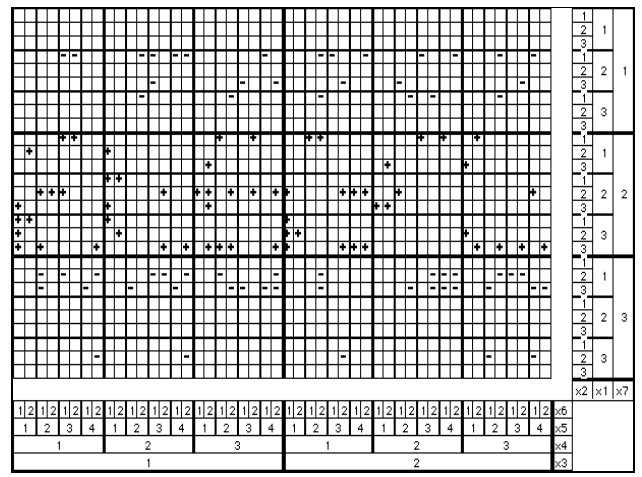
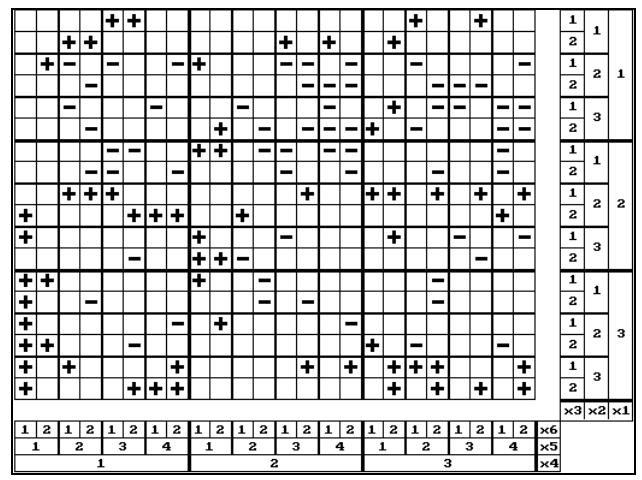


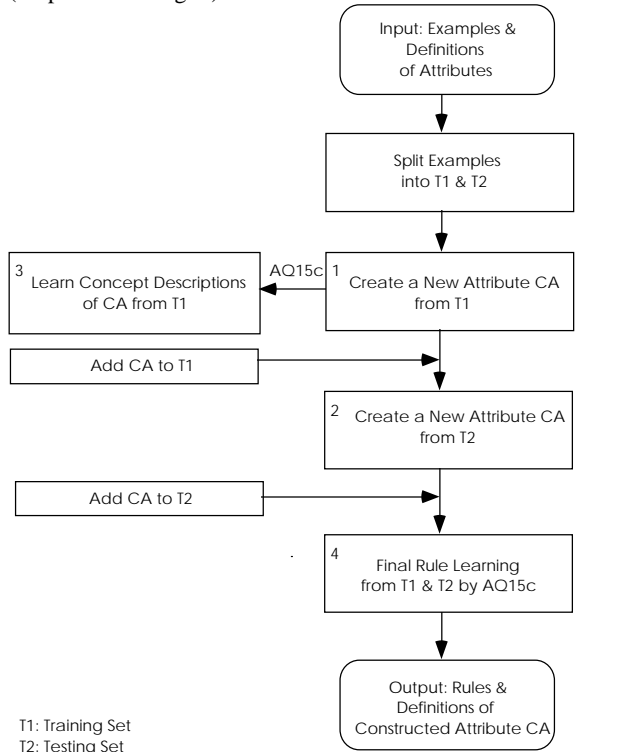
Figure 2. Diagrammatic visualization of the first MONK problem representation space: A) the initial representation space; B) the improved representation space due to the AqBC approach

**AutoClass: An Unsupervised Bayesian Classifier**  
 AutoClass is an unsupervised Bayesian classification system that looks for a maximum posterior probability classification. The system infers classes based on Bayesian statistics, deriving a belief network via probability theory.

The idea of accumulating and applying auxiliary evidence here can be mapped into the constructive induction mechanism that employs a new attribute which summarizes the data patterns. The new attribute's degree of belief is very high because it is generated from the best model of Bayesian classification. Therefore, this new attribute can potentially reorganize and improve the knowledge representation space. The theory of how Bayesian learning is applied in AutoClass is described in [2], [5].

**The AqBC Multistrategy Learning Methodology**

AqBC is a multistrategy knowledge discovery approach that combines unsupervised Bayesian classification with supervised inductive learning. Figure 3 shows the general structure of the AqBC approach. AqBC can be applied to two different goals. First, AutoClass provides classifications that help the user generate "expert knowledge" for a potential target concept when a data set without predefined classifications is given to the supervised inductive learning system. The inductive learning system AQ15c learns the concept descriptions of these classifications (step3 in Fig. 3). Second, when the data set is already divided into classes, AqBC repeatedly searches for the best classifications using AutoClass, then uses the best one for creating and modifying more suitable and meaningful knowledge representation spaces. Figure 4 shows the constructive induction process in which a new attribute is created for training and testing. In the first phase, a new attribute, in which the values are class labels generated by AutoClass, is added to the data table under the name of CA (Step 1 & 2 in Fig. 3).



T1: Training Set  
 T2: Testing Set  
 CA: New attribute with values that are classes generated by AutoClass  
 In step 1, concept attribute (C) is included in T1. In step 2, C is not included in T2.

Figure 3. General structure of AqBC approach

Note that one of the attributes provided to AutoClass is the original target concept  $C$  (Fig. 4-a). However,  $C$  is not included when we are dealing with the testing data set (Fig. 4-b) since the system is not supposed to know the given class.

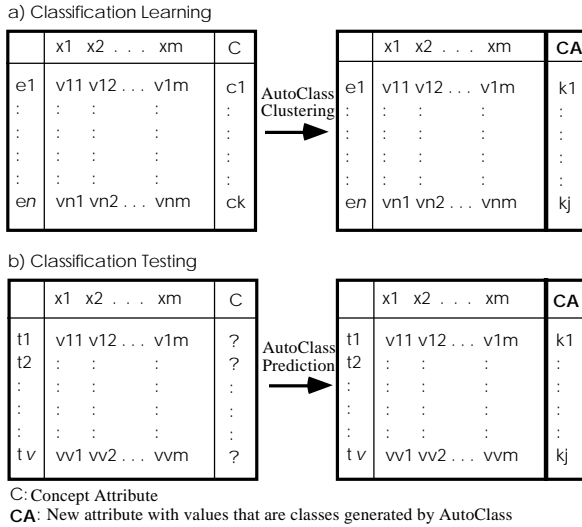


Figure 4. Creating a new attribute for a) training and b) testing data set by use of AutoClass, where,

- vnm : a value of attribute xm of nth example
- $c1 \leq C \leq ck$  : C is a Target Concept Attribute.
- k is the # of target concepts.
- $k1 \leq CA \leq kj$ : CA is a classification label generated by AutoClass
- j is the # of classifications.

The modified data table generated by this phase, which now includes CA, is used by AQ15c to learn the new concept description for C. In other words, the concept description of the original classification is learned from the modified knowledge representation space. The only change is the augmentation of the table with the new attribute. It is at this time that a separate AQ15c run generates the concept descriptions for the classes represented by CA.

In the second phase, the testing data set, an unseen data set with the original concept attribute C containing the value “?” (don’t care) for all examples, is input to the AutoClass classifier trained in the first phase (see Figure 4b).

By not providing the original classifications, the original properties of the testing data set can be preserved and AutoClass’s ability to correctly classify previously unseen examples can be verified. The technique of constructing classification labels is based on previous constructive induction methodology [14].

The improvement of the knowledge representation space is already demonstrated with the MONK1 problem in the previous section. Such an approach can potentially allow us to address large, complicated data sets not recognized well by a supervised inductive learner alone, by having an unsupervised Bayesian classifier as an oracle. Then the robust supervised inductive learner learns the target concept descriptions, providing clear, understandable rules describing these classes. Thus, the system takes advantage of the strengths of its component subsystems.

#### 4. AN EXPERIMENTAL APPLICATION TO THE US CENSUS DATA

##### US Census Data

The demographics database (Table 1) is adapted from US census data on all US cities with a 1990 population of over 200,000 residents. 77 records containing the above 36 attributes were used for the experiments.

The sizes of the domains of these attributes vary widely. Some are binary, such as the last four attributes; others have as many as 70 different linear values. Most of the attributes have 10-20 values. The dataset thus introduces a wide variety of problems for any system that must deal with it, particularly in view of the small number of examples. Given such a dataset, an important problem is how to organize it in such a way so that useful knowledge can be extracted from it.

##### Learning Concept Descriptions from AutoClass classifications

AQ15c learns the concept description of the classifications obtained from the AutoClass (attribute CA), as described earlier. The following rules describe the two values of the CA attribute, representing the two concepts discovered by AutoClass. Note that, for all rules presented below, the values have been mapped back to the original data from the normalized values presented to the system.

Table 1. Definitions of attributes

1. % Asian Residents	9. Population	17. % Jobs held are in manufacturing	23. % Dwellings that are condominiums	30. % of children in elementary or secondary school
2. % Black	10. Infant Mortality	18. Change in labor force from 1980	24. Home value	31. % Adults over 55
3. Population Density	11. Crime Rate	19. % Females in workplace	25. % Population who rent	32. % Holding bachelors degrees
4. % Elderly Residents	12. % Single Dwellers	20. % Commuters using public transit	26. Poverty rate	33. Has major league baseball team
5. % Foreign Residents	13. % Single Parents	21. Unemployment %	27. % Housing units built before 1939	34. Has NFL football team
6. % Residents speaking foreign language	14. % Hispanic Residents	22. Rent	28. Income	35. Has NBA basketball team
7. Population Growth	15. July Temperature (median)		29. % Residents on public assistance	36. Is a state capital
8. Land Area	16. Precipitation (annual, inches)			

#### class0-outhypo

- 1 [black = 4..76%] & [elderly = 4..18%] & [pop\_growth < 32%] & [population > 385,000] (t:41, u:26)
- 2 [manufacturing\_jobs = 7..27%] & [rent < \$475/mo] & [poverty = 9..17%] (t:20, u:5)

**class0** appears to represent cities with moderate to large populations and without explosive growth OR with low average monthly rent and moderate poverty rates. When we look at the actual cities that fall into this class, we see most are cities that are very large, with stable or declining populations – particularly those of the Eastern seaboard, the old South and the Rust Belt. Of the 32 largest cities in the US, 30 fall into this class. The only exceptions are the high tech, rapidly growing San Jose, and El Paso, which is also growing extremely rapidly, due primarily to Hispanic immigration.

#### class1-outhypo

- 1 [foreign\_speak = 4..70%] & [land\_area < 242 sq miles] & [renters = 0..21%] & [old\_housing = 0..19%] & [income = \$22K..46K] (t:22, u:19)
- 2 [population = 326K..385K] & [infant\_mortality = 6..18] & [renters = 6..40%] & [state\_capital = NO] (t:11, u:8)
- 3 [black = 43%] & [foreign < 19%] (t:1, u:1)

**class1** appears to represent small cities with relatively low rental rates or cities with moderately small populations that aren't state capitals. Most of these cities are either in the Sun Belt or are small cities outside of the old South and Rust Belt.

#### Using AqBC for Knowledge Discovery

AQ15c learns the abstract concept descriptions of the given classifications from AutoClass in the first phase. Now we augment the original knowledge representation space with the “class” label attribute, allowing AQ to learn new knowledge which was difficult to extract from the original representation space. We can now choose our target concepts from among any of the system attributes. In this case, we choose two concepts based on the *population* attribute, where these concepts, *large cities* and *moderately sized cities*, are represented by two classes. The first class is population over 385,000 and the second class is population from 200,000 to 385,000.

The following experiments use the constructed attribute from the two class clustering.

#### population\_0-outhypo

- 1 [black=7..75%] & [foreign = 2..59%] & [home\_value = 2..32] & [class=0] (t:34, u:24)
- 2 [pop\_density > 2250] & [single = 24..31%] & [july\_temperature = 74..93F] & [precipitation < 56] & [manufacturing\_jobs < 24%] & [change\_in\_labor < 43%] & [renters < 56%] (t:14, u:9)

- 3 [foreign = 4..27%] & [foreign\_speak > 1%] & [pop\_growth < 32%] & [july\_temperature = 0..24F] & [precipitation = 4..43] & [change\_in\_labor > -7%] & [renters < 56%] (t:11, u:4)
- 4 [single\_parent = 17%] (t:1, u:1)

#### population\_1-outhypo

- 1 [hispanic > 66%] & [manufacturing\_jobs = 6..26%] & [class=1] (t:21, u:15)
- 2 [pop\_density < 6750] & [land\_area < 129] & [manufacturing\_jobs < 22%] & [unemployment ≠ very high] & [renters = 41..61%] & [holds\_bachelors < 28%] & [has\_nba\_team = NO] (t:9, u:5)
- 3 [single = 20..23%] (t:5, u:1)

The constructed attribute separates large cities from small ones almost perfectly. As stated above, 30 of the 32 largest cities (and 34 of the 38 largest cities that fit into population class 0) are labeled 0 for the constructed attribute “class”. This greatly simplifies the rule learner’s task for discovering simple rules representing the target concept “population”.

The Census database has a much larger attribute set than the MONK1 problem, and its attributes also have a wider range of values. Because of this, the Census database provides a somewhat more realistic case study. More than half the attributes of the database are dropped altogether in the final rules, while still capturing the target concept. In addition, the constructed attribute plays a major role in discriminating the two classes; the **class** attribute is part of the most powerful conjunction in each class description. Thus, constructive induction turns out to be important to understanding the target concept.

## 5. CONCLUSIONS

This paper presents a multistrategy approach for knowledge discovery that integrates supervised AQ inductive rule learning with unsupervised Bayesian classification for creating a meaningful knowledge representation space and discovering high quality knowledge. The patterns acquired by the Bayesian classifier make little sense until AQ provides an understandable descriptive representation that encompasses new knowledge. In other words, the data and information presented to AqBC becomes knowledge when the patterns acquired by the system are realized into understandable concepts. This new knowledge can then be applied in organizational decision making.

The AqBC approach uses the methodology of constructive induction, which modifies and improves the representation space with the aid of classifications obtained from the unsupervised Bayesian classifier, AutoClass. This methodology was applied to the MONK1 problem and a US census dataset. In the experiments, AqBC constructed new attributes used in relatively simple, meaningful rules describing the target concepts. The problems inherent in the nonsymbolic nature of both the Bayesian and the distance based statistical techniques

are often mitigated by applying AQ to learn descriptive representations for both the resulting clusters and the target concepts. This multistrategy approach produces important new knowledge organization as it creates human-understandable class descriptions and new attributes that simplify those class descriptions.

A key result of this research is the generation of a simple, descriptive taxonomy for a database that can be used for sensibly partitioning into smaller databases. Another potential experiment would be to use the taxonomy produced via the AqBC method to perform supervised classification with many different decision attributes from the current set rather than just the population attribute. Other future research will focus on developing new strategies combining various statistical and conceptual [10] classification methods for constructing better knowledge representation spaces from large data sets.

Not only Bayesian classifiers benefit from the use of AQ15c to learn the descriptive representations of the generated classes. In addition to the distance-based clustering methods like K-means centroid & Ward hierarchical clustering, other subsymbolic systems such as SOFMs [6] and k-NN methods can be employed as the unsupervised classification engine. Multiple engines may also be used to perform additional constructive induction prior to unsupervised classification. For example, AQ17-DCI [1] and AQ17-HCI [14] construct new features based on interrelationships among existing ones.

Varying the rule learner based on the application may also prove productive. If mathematical formulae are desired instead of conjunctive rules, a system such as ABACUS [3] could be employed in place of AQ15c. There are still many challenging real world applications to which this multistrategy approach could be applied for new knowledge discovery.

## 6. ACKNOWLEDGMENTS

The authors thank Dr. Ryszard Michalski for his comments and the AutoClass group for making their software readily available. The first author greatly appreciates the support of a Doctoral Fellowship from the School of Information Technology and Engineering at George Mason University. This research was conducted in the Machine Learning and Inference Laboratory at George Mason University. The Laboratory's research activities are supported in part by the National Science Foundation under grants IRI-9020266 and DMI-9496192, in part by the Advanced Research Projects Agency under grants F49620-92-J-0549 and F49620-95-1-0462, administered by the Air Force Office of Scientific Research, and in part by the Office of Naval Research under grant N00014-91-J-1351.

## 7. REFERENCES

[1] E. Bloedorn and R. S. Michalski, "Constructive Induction from Data in AQ17-DCI: Further Experiments," *Reports of the Machine Learning and Inference Laboratory*, MLI 91-12, George Mason University, Fairfax, VA, 1991.

[2] P. Cheeseman and J. Stutz, "Bayesian Classification (AutoClass): Theory and Results," In *Advances in Knowledge*

*Discovery and Data Mining*. U. M. Fayyad, G. Piatetsky-Shapiro, P. Smyth, and R. Uthurusamy, eds., AAAI Press, Menlo Park, 1996.

[3] B. C. Falkenhainer and R. S. Michalski, "Integration Quantitative and Qualitative Discovery in the ABACUS System," In *Machine Learning: An Artificial Intelligence Approach, Vol. III*. Y. Kodratoff and R. S. Michalski, eds., Morgan Kaufmann, San Mateo, 1990.

[4] U. M. Fayyad, G. Piatetsky-Shapiro, P. Smyth, and R. Uthurusamy, eds. *Advances in Knowledge Discovery and Data Mining*, The AAAI Press, Menlo Park. 1996.

[5] R. Hanson, J. Stutz, and P. Cheeseman, "Bayesian Classification Theory," Technical Report FIA-90-12-7-01, NASA Ames Research Center, Artificial Intelligence Branch, 1991.

[6] T. Kohonen, "*Self Organizing Maps*," Springer-Verlag, Heidelberg, 1995.

[7] S. W. Lee, "Multistrategy Learning: An Empirical Study with AQ + Bayesian Approach," *Reports of the Machine Learning and Inference Laboratory*, MLI 96-10, George Mason University, Fairfax, VA, 1996.

[8] R. S. Michalski, "Pattern Recognition as Knowledge-Guided Computer Induction," Technical Report No. 927, Department of Computer Science, University of Illinois, Urbana-Champaign, 1978.

[9] R. S. Michalski, "A Theory and Methodology of Inductive Learning," *Artificial Intelligence*. Vol.20, 111-116, 1983a.

[10] R. S. Michalski and R. E. Stepp, "Automated Construction of Classifications: Conceptual Clustering Versus Numerical Taxonomy," *IEEE Transactions on PAMI* 5:4, 396-410, 1983b.

[11] R. S. Michalski, I. Mozetic, J. Hong, and N. Lavrac, "The Multi-Purpose Incremental Learning System AQ15 and its Testing Application to Three Medical Domains," In *Proceedings of AAAI-86*, 1041-1045, Philadelphia, PA, 1986.

[12] S. B. Thrun, et al. "The MONK's problems: A Performance Comparison of Different Learning Algorithms," Carnegie Mellon University, 1991.

[13] J. Wnek, K. Kaufman, E. Bloedorn, and R. S. Michalski, "Inductive Learning System AQ15c: The Method and User's Guide," *Reports of the Machine Learning and Inference Laboratory*, MLI 95-4, George Mason University, VA, 1995.

[14] J. Wnek and R. S. Michalski, "Hypothesis-driven Constructive Induction in AQ17-HCI: A method and experiments," *Machine Learning*, Vol.14, No.2, 139-168, 1994.

[15] J. Wnek, "DIAV 2.0 User Manual: Specification and Guide through the Diagrammatic Visualization System," *Reports of the Machine Learning and Inference Laboratory*, MLI 95-4, George Mason University, Fairfax, VA, 1995.