

An Adjustable Description Quality Measure for Pattern Discovery in Large Databases Using the AQ Methodology*

KENNETH A. KAUFMAN

Machine Learning and Inference Laboratory, George Mason University,
Fairfax, VA 22030

kaufman@mli.gmu.edu

RYSZARD S. MICHALSKI

Machine Learning and Inference Laboratory, George Mason University,
Fairfax, VA, and
Institute of Computer Science, Polish Academy of Sciences, Warsaw, Poland

michalsk@mli.gmu.edu

Abstract. In concept learning and data mining tasks, the learner is typically faced with a choice of many possible hypotheses or patterns characterizing the input data. If one can assume that training data contain no noise, then the primary conditions a hypothesis must satisfy are consistency and completeness with regard to the data. In real-world applications, however, data are often noisy, and the insistence on the full completeness and consistency of the hypothesis is no longer valid. In such situations, the problem is to determine a hypothesis that represents the best trade-off between completeness and consistency. This paper presents an approach to this problem in which a learner seeks rules optimizing a *rule quality criterion* that combines the rule coverage (a measure of completeness) and training accuracy (a measure of inconsistency). These factors are combined into a single rule quality measure through a *lexicographical evaluation functional* (LEF). The method has been implemented in the AQ18 learning system for natural induction and pattern discovery, and compared with several other methods. Experiments have shown that the proposed method can be easily tailored to different problems and can simulate different rule learners by modifying the parameter of the rule quality criterion.

Keywords: Machine Learning, Data Mining, Learning from Noisy Data, Natural Induction, AQ Learning, Decision Rules, Separate and Conquer

1. Introduction

In concept learning and data mining tasks, a typical objective is to determine a general hypothesis characterizing the training instances, which will classify future instances as correctly as possible. For any non-trivial generalization problem, one can usually generate a large number of plausible hypotheses that are complete and consistent with regard to the training data (i.e., covering all positive examples and no negative examples). If one can assume that training data contain no noise, hypothesis consistency and completeness are justifiable preconditions for admissibility of a hypothesis. In real-world applications, however, data are often noisy (contain errors) and/or inconsistent (contain contradictions); therefore, the completeness and consistency criteria are no longer apply. In such situations, one may seek a hypothesis that is maximally consistent at the expense of completeness, maximally complete at the expense of inconsistency, or representing some combination of the two criteria.

* This research was conducted in the Machine Learning and Inference Laboratory at George Mason University. The Laboratory's research activities have been supported in part by the National Science Foundation under Grants No. IIS-9904078 and IRI-9510644, in part by the Defense Advanced Research Projects Agency under Grant No. F49620-95-1-0462 administered by the Air Force Office of Scientific Research, and in part by the Office of Naval Research under Grant No. N00014-91-J-1351.

The problem then arises as to what kind of trade-off should be chosen between completeness and consistency, or, more generally, what criteria should be used to guide the learner in problems with noisy and/or inconsistent data?

To illustrate this problem, suppose that a training set contains 1000 positive (P) and 1000 negative (N) examples of the concept to be learned. Suppose further that the system generated two hypotheses, one covering 600 positive (p) and 2 negative (n) examples, and another one covering 950 positive and 20 negative examples. Which is better?

Clearly, the choice of which hypothesis to select is not domain-independent. In some domains, the first hypothesis, with a 60% completeness level (p / P) and 99.7% consistency level ($p / (p + n)$), may be preferred because it represents a more consistent pattern (henceforth, instead of “completeness level” and “consistency level” will use simply the terms “consistency” and “completeness” whenever it leads to no confusion). In other domains, the second hypothesis, with 95% completeness and 98% consistency, may be viewed as better, because it represents a more dominant pattern. It is also possible that a third principle such as which rule is most easily articulated to or understood by the user will sway the choice.

This paper explores the issues related to the trade-off between hypothesis and consistency in the case of noisy or inconsistent training data, and proposes a single description quality measure that reflects this trade-off. The learning process is presented as a search for a hypothesis that maximizes the description quality measure. To get a better understanding of the proposed measure, rankings of hypotheses by this measure and other known criteria have been calculated and compared. The concluding sections discuss the implementation of the proposed method in the AQ18 rule learning system for natural induction and pattern discovery (Michalski, 1999).

2. How to Choose the Best Hypothesis

In the progressive covering approach to concept learning (also known as separate-and-conquer), if the training data can be assumed to be noise-free, then primary conditions for admitting a hypothesis are consistency and completeness with regard to data. Other factors, such as computational simplicity, description comprehensibility, and focus on preferred attributes, are considered after the consistency and completeness criterion is satisfied. If the training data contain errors (class errors or value errors) or inconsistency (the same example occurs in more than one training class), some degree of inconsistency and incompleteness of the rulesets can be not only acceptable, but also desirable (e.g., Bergadano et al, 1992). In such situations, a selection hypothesis criterion is typically some function of the number of positive and negative examples covered by a rule.

For example, in the RIPPER program (Cohen, 1995), the criterion is to maximize:

$$(p - n) / (P + N) \tag{1}$$

where p and n are the numbers of positive and negative examples covered by the rule, and P and N are the numbers of positive and negative examples in the entire training set, respectively. In Section 5, we will explore RIPPER's and several other criteria used by the machine learning and data mining communities, and compare their rankings of various hypotheses.

Although such criteria present heuristics for choosing among rules, it is unrealistic to assume that any single criterion will fit all practical situations. For different problems, different quality criteria of a hypothesis (a

description, a rule, etc.) may lead to the best performance. The problem of determining the best hypothesis can be characterized generally as a problem of optimizing an assumed a priori measure of hypothesis quality.

Such a measure (criterion) may combine various constituent criteria such as completeness (here, the percentage of positive examples covered by it), and the consistency, or a reciprocal measure, inconsistency (here, the percentage of positive or negative examples covered, respectively). Several specific criteria combining completeness and consistency have been described in the literature (e.g., Bruha, 1997). It is, however, unrealistic to expect that any such criterion will fit all practical problems. For different problems, different criteria of description quality may lead to better performance.

In this paper, learning is viewed as a process of optimizing a description quality criterion that best reflects the characteristics of the problem at hand. This view has been implemented in the learning system AQ18, which allows a user to apply a combination of different *elementary criteria*, each representing one aspect of the hypothesis (ruleset) being learned. These elementary criteria are selected from a menu of available elementary preference criteria. The selected criteria are invoked sequentially; when one criterion is insufficient for distinguishing between two or more candidate rules, the next one is used, etc., until the rule of the highest rank is determined.

This process is described precisely by the *lexicographical evaluation functional* (LEF; Michalski, 1983), defined as a sequence:

$$\langle (c_1, \tau_1), (c_2, \tau_2), \dots, (c_n, \tau_n) \rangle \quad (2)$$

where c_i represents the i th elementary criterion, and τ_i is the *tolerance* associated with c_i . The latter defines the range (either absolute or relative) within which a candidate rule's c_i evaluation value can deviate from the best evaluation value of this criterion in the current set of rules (generally, descriptions). Let us assume, for example, that we have a set of hypotheses, S , and that there are just two elementary criteria, one, to maximize the completeness (or coverage), and second, to minimize inconsistency. Let us assume further that hypotheses with coverage within 10% of the maximum coverage achievable by any single rule in S is acceptable, and that if two or more hypotheses satisfy this criterion, the one with the lowest inconsistency is to be selected. The above rule selection process can be specified by the following LEF:

$$\text{LEF} = \langle (\text{coverage}, 10\%), (\text{inconsistency}, 0\%) \rangle \quad (3)$$

It is possible that after applying both criteria, more than one hypothesis remains in the set of candidates. In this case the one that maximizes the first criterion is selected.

The advantages of the LEF approach are that it is very simple to apply and very efficient, so that it can be effectively applied in the case of a very large number of candidate hypotheses. In (3), coverage and inconsistency are used sequentially, and according to the specified tolerance levels. An alternative approach is to combine these two criteria into a single numerical measure, as is done in various machine learning programs. To get an insight into how these programs rank candidate hypotheses, Section 5 presents experimental results from applying them to selected hypothetical problems.

3. Completeness, Consistency and Consistency Gain

As mentioned above, in data sets that may contain noise, which typically occurs in real-world data mining applications, full consistency and completeness are not necessary. Full consistency and completeness can also be

ignored when one seeks strong patterns in the data and allows for exceptions. In such cases, one seeks descriptions that optimize some quality criterion. This paper describes an approach to learning such descriptions implemented in the AQ18 system for natural induction and pattern discovery. This system applies a progressive covering approach, also known as separate-and-conquer, in which partial data descriptions (in the form of rules) are determined sequentially. This approach thus requires an evaluation of single rules at each step of the process.

As the main purpose of the learned rules is to be able to classify correctly future, unknown cases, a useful measure of rule quality is the *testing accuracy*, that is, the accuracy of classifying testing examples, which are different from the training examples. During a learning process, testing examples by definition are not being used, therefore one needs a criterion that approximates well the testing accuracy on the basis of only training examples. Before we propose such a measure, we need to introduce some notation and terminology.

Let P and N denote the total number of positive and negative examples, respectively, of some concept or decision class in a training set. Let R be a rule or a set of rules (a ruleset) generated to cover examples of that class, and p and n be the number of positive and negative examples covered by R , respectively. For the given rule, the ratio p / P , denoted $\text{compl}(R)$, is called the *completeness* or *relative coverage* of R . The ratio $p / (p + n)$, denoted $\text{cons}(R)$, is called the *consistency* or *training accuracy* of R , and $n / (p + n)$, denoted $\text{inc}(R)$, the *inconsistency* or *training error rate*. If the completeness of a ruleset (a set of rules for a single class) is 100%, then the ruleset is a *complete cover* of the training examples. If the inconsistency of the ruleset is 0%, then it is a *consistent cover*.

We have tried to maintain a terminology that is consistent with both the literature and intuitive understanding. With issues such as these of interest in different research communities, there is no agreement on the notation and terminology. In addition to the terminology used in this report, the terminologies seen in Bruha (1997) and Fayyad et al (1996) are representative of sizable segments of the machine learning and data mining communities respectively.

Let us return to the question posed in the Introduction as to which is preferable: a rule with 60% completeness and 99.7% consistency, or a rule with 95% completeness and 98% consistency. As indicated earlier, the answer depends on the problem at hand. In some application domains, notably in science, a rule (law) must be consistent with all the data, unless some of the data are found erroneous. In other applications, in particular, data mining, one may seek strong patterns that hold frequently, but not always. Therefore, there is no single measure of rule quality that would be good for all problems. Instead, we seek a flexible measure that can be easily changed to fit any given problem at hand.

In general, one may assume that which rule is “better” depends, all other factors being equal, on a function of the completeness and the consistency. A third factor, rule simplicity, is also important, especially in cases in which two rules rank similarly based on completeness and consistency. The simplicity can be taken into consideration by properly defining the LEF criterion.

How then can we define such a measure of rule quality? One approach to quantifying such considerations is the *information gain* criterion used in decision tree learning for selecting attributes (e.g., Quinlan, 1986). This criterion can also be used for selecting a decision rule, as such a rule can be viewed as a binary attribute whose value represents whether or not an instance is in the portion of the event space covered by the rule. Given a set

of events divided into positive, P , and negative, N , event sets in an event space E , the entropy, or alternatively, the expected information from a message defining the class of an event is defined as:

$$\text{Info}(E) = -((P / (P + N)) \log_2(P / (P + N)) + (N / (P + N)) \log_2(N / (P + N))) \quad (4)$$

The expected information of using rule R to partition the event space into regions covered and not covered by the rule is defined as

$$\text{Info}_R(E) = ((p + n) / (P + N)) \text{Info}(R) + ((P + N - p - n) / (P + N)) \text{Info}(\sim R) \quad (5)$$

Then the information gained by using rule R to partition the event space is:

$$\text{Gain}(R) = \text{Info}(E) - \text{Info}_R(E) \quad (6)$$

This measure is, as desired, a function of completeness and consistency; the higher a rule's completeness or consistency, the better it will score. Information gain has, however, one major disadvantage as a rule evaluator. It relies not only on the informativeness of the rule, but also the informativeness of the complement of the rule. That is, it takes into consideration the entire partition created by the rule, rather than just the space covered by it. This concern is especially valid if there are more than two decision classes. In such a situation, a rule may be highly valuable for classifying examples of one specific class, even if it does little to reduce the entropy of the training examples in other classes.

As an example, consider the problem of distinguishing upper-case letters of the English alphabet. In this case, the rule, "If a capital letter has a tail, it is the letter Q" is simple, with a perfect or near-perfect coverage and consistency for the Q class. As it is a very specific rule, tailored toward one class, the above gain measure applied to it will produce a low score. Another limitation of the information gain measure is that it does not provide the means for modifying it to fit different problems which may require different relative weights to be placed on consistency and completeness.

Before proposing another measure, let us observe that relative frequency of positive and negative examples in the training set of a given class should also be a factor in evaluating a rule. Clearly, a rule with 15% completeness and 75% consistency could be quite attractive if the total number of positive examples was very small, and the total number of negative examples was very large. On the other hand, the same rule would be uninformative if P was very large and N was very small.

The distribution of positive and negative examples in the whole training set can be measured by the ratio $P / (P + N)$. The distribution of positive and negative examples in the set covered by the rule can be measured by the consistency $p / (p + n)$. Thus, the difference between these values, $(p / (p + n)) - (P / (P + N))$, which reflects the gain of the consistency over the dataset's distribution as a whole, can be normalized by dividing it by $(1 - (P / (P + N)))$, or equivalently $N / (P + N)$, so that in the case of identical distribution of positive and negative events in the set covered by the rule and in the training set, it returns 0, and in the case of perfect training accuracy, it will return 1. This normalized consistency measure thus shares the *independence property* with statistical rule quality measures (Bruha, 1997).

This measure thus provides an indication of the benefit of the rule, based on its consistency, over making a random guess, and allows for the possibility of negative values, in accordance with our assertion that a rule less accurate than the random guess based the example distribution has a negative benefit. Reorganizing the normalization term, we define the *consistency gain* of a rule R , $\text{consig}(R)$, as:

$$\text{consig}(\mathbf{R}) = ((p / (p + n)) - (P / (P + N))) * (P + N) / N \quad (7)$$

4. A Formula for Rule Quality

In developing a rule quality measure, one may assume the desirability of maximizing both the completeness, $\text{compl}(\mathbf{R})$, and the consistency gain, $\text{consig}(\mathbf{R})$. Clearly, a rule with a higher $\text{compl}(\mathbf{R})$ and a higher $\text{consig}(\mathbf{R})$ is more desirable than a rule with lower measures. A rule with either $\text{compl}(\mathbf{R})$ or $\text{consig}(\mathbf{R})$ equal to 0 is worthless. It makes sense, therefore, to define a rule quality measure that evaluates to 1 when both of these components reach maximum (have value 1), and 0 when either is equal to 0.

A simple way to achieve such a behavior is to define rule quality as a product of $\text{compl}(\mathbf{R})$ and $\text{consig}(\mathbf{R})$. Such a formula, however, does not allow one to weigh these factors differently in different applications. To achieve this flexibility, we introduce a weight, w , defined as the percentage of the description quality measure to be borne by the completeness condition. Thus, the final form of the *description quality*, $Q(\mathbf{R}, w)$ with weight w , or just $Q(w)$, if the rule \mathbf{R} is implied, is:

$$Q(\mathbf{R}, w) = \text{compl}(\mathbf{R})^w * \text{consig}(\mathbf{R})^{(1-w)} \quad (8)$$

By changing parameter w , one can change the relative importance of the completeness and the gain in consistency to fit a given problem. It can be seen that when $w < 1$, $Q(w)$ satisfies the constraints listed by Piatetsky-Shapiro (1991) regarding the behavior of rule evaluation criteria:

1. The rule quality should be 0 if the example distribution in the space covered by the rule is the same as in the entire data set. Note that $Q(\mathbf{R}, w) = 0$ when $p/(p+n) = P/(P+N)$, assuming $w < 1$.
2. All other things being equal, an increase in the rule's coverage should increase the quality of the rule. Note that $Q(\mathbf{R}, w)$ increases monotonically with p .
3. All other things being equal, the quality of the rule should decrease when the ratio of covered positive examples in the data to either covered negative examples or total positive examples decreases. Note that $Q(\mathbf{R}, w)$ decreases monotonically as either n or $(P - p)$ increases, when $P + N$ and p remain constant.

The formula cited as the simplest one satisfying the above three criteria in the notation used by Piatetsky-Shapiro is none other than $\text{consig}(\mathbf{R})$ without the normalization factor, and multiplied by $(p + n)$ (Piatetsky-Shapiro, 1991). The advantage of incorporating the component of $\text{compl}(\mathbf{R})$ in (8) is that it promotes rules with high coverage, and by that, rules that are applicable to a larger number of cases. The next section compares the proposed $Q(w)$ rule evaluation method with other methods, and Sections 6 and 7 discuss its implementation in AQ18.

5. Empirical Comparison of Evaluation Methods

In order to develop a sense of how the $Q(w)$ rule rankings compare to those done by other methods used in machine learning systems, we performed a series of experiments on different datasets. In the experiments we used the $Q(w)$ method with different weights, the information gain criterion (Section 3), the PROMISE method (Baim, 1982; Kaufman, 1997), and the methods employed in CN2 (Clark and Niblett, 1989), IREP (Fürnkranz and Widmer, 1994) and RIPPER (Cohen, 1995) rule learning programs. In describing these methods, we will use the same notation for positive and negative examples as was used in the earlier sections of this paper.

As was mentioned above, the information gain criterion is based on the entropy of the examples in the area covered by a rule, the area not covered by the rule, and the event space as a whole. Like the information gain criterion, the PROMISE method (Baim, 1982; Kaufman, 1997) was developed to evaluate the quality of attributes. It can be used, however, for rule evaluation by considering a rule to be a binary attribute that splits the space into the part covered by the rule and the part not covered by it. The application of PROMISE to such an attribute reduces to the description given below in which:

$$M_+ = \max(p, n)$$

$$M_- = \max(P - p, N - n)$$

$$T_+ = P \text{ if } p > n, N \text{ if } p < n, \text{ and } \min(P, N) \text{ if } p = n$$

$$T_- = P \text{ if } P - p > N - n, N \text{ if } P - p < N - n, \text{ and } \min(P, N) \text{ if } P - p = N - n$$

PROMISE will return a value of $(M_+ / T_+) + (M_- / T_-) - 1$ (the last term is a normalization factor to make the range 0 to 1). It should be noted that when M_+ and M_- are based on the same class (for example, the positive class, as is the case when $p > n$ and $P - p > N - n$), PROMISE will return a value of zero. Hence, it is not a useful measure of rule quality in domains in which the positive examples significantly outnumber the negative ones. Note also that when $P = N$ and p exceeds n (the latter presumably occurs in any rule of value in an evenly distributed domain), the PROMISE value reduces to:

$$(p - n) / P \tag{9}$$

To see this, note that when $P = N$, $(p / P) + ((N - n) / N) - 1$ can be transformed into $(p / P) + ((P - n) / P) - 1$, which is equivalent to (9).

CN2 (Clark and Niblett, 1989) builds rules using a beam search, as does the AQ-type learner, on which it was partially based. It evaluates rules based on an entropy measure, as it attempts to minimize, in the case of two decision classes:

$$-((p / (p + n)) \log_2(p / (p + n)) + (n / (p + n)) \log_2(n / (p + n))) \tag{10}$$

This expression involves only the consistency, $p / (p + n)$; it does not involve any completeness component. Thus, a rule that covers 50 positive and 5 negative examples is deemed of identical value to another rule that covers 50,000 positive and 5000 negative examples. Although (10) has a somewhat different form than the rule utility part of $Q(w)$, CN2's rule evaluation can be expected to be similar to $Q(0)$ (utility only). Indeed, in the examples shown below, the two methods provide identical rule rankings.

If there are more than two decision classes, the entropy terms are summed. Nonetheless, the above comments regarding no consideration of rule completeness hold true.

A later version of CN2 (Clark and Boswell, 1991) offered a new rule quality formula based on a Laplace error estimate. This formula is closely tied to a rule's consistency level, while completeness still plays a minimal role.

IREP's formula for rule evaluation (Fürnkranz and Widmer, 1994) is:

$$(p + N - n) / (P + N) \tag{11}$$

RIPPER, as was mentioned in Section 2, uses a slight modification of the above formula:

$$(p - n) / (P + N) \tag{12}$$

Note that RIPPER’s evaluation will not change when P changes, but $P + N$ stays constant. In other words, its scores are independent of the distribution of positive and negative examples in the event space as a whole. While this therefore evaluates a rule on its own merits, the evaluation does not factor in the benefit provided by the rule based on the overall distribution of classes.

Furthermore, since P and N are constant for a given problem, a rule deemed preferable by IREP will also be preferred by RIPPER. Thus, these two measures produce exactly the same ranking; in comparing different measures, we therefore only show RIPPER’s rankings below. Comparing (12) to (9), one can notice that RIPPER evaluation function returns a value equal to half of the PROMISE value when $P = N$ and p exceeds n . Thus, in such cases, the RIPPER ranking is the same as the PROMISE ranking.

We compared the above methods on three datasets, each consisting of 1000 training examples. Dataset A has 20% positive and 80% negative examples, Dataset B has 50% positive and negative examples, and Dataset C has 80% positive examples and 20% of negative examples. In each dataset, rules with different coverage and training accuracy were ranked using the following criteria: Information Gain, PROMISE, RIPPER, CN2 (the initial implementation), $Q(0)$, $Q(.25)$, $Q(.5)$, $Q(.75)$, and $Q(1)$.

Results are summarized in Table 1. The leftmost column identifies the dataset, the next two give the numbers of positive and negative examples respectively covered by a hypothetical rule, and the remaining columns list the evaluations and ranks on the dataset of the rules by the various methods. Most of the values are normalized into a 0-1 range, although as was discussed in Section 4, a Q value could fall below 0 if the rule gave support to the negative class; Information Gain may also not fall into such a range. We reiterate, however, that the ranking of rules is more significant than their particular quality values.

There is, of course, no one answer regarding which ranking is superior. It should be noted, however, that by modifying the Q weights, one can tailor the rule evaluation criterion according to the problem at hand.

6. Methods for Introducing $Q(w)$ in AQ18

There are four points during the rule generation process at which AQ18 evaluates and selects candidate rules. They are:

1. *Star generation.* During this phase, AQ uses a beam search strategy to find near-optimal generalizations of a seed example. This has implied near-optimality in the context of full consistency (barring ambiguity handling).
2. *Star completion and rule selection.* After a set of consistent rules has been built through iterations of the star generation step, the best one is selected and added to the preliminary output ruleset.
3. *Rule trimming.* Once the best rule that covers the seed example has been selected, the rule may be specialized to reflect the user’s output preferences. The final rule is selected and added to the preliminary final ruleset.
4. *Rule truncation.* Once all of the rulesets have been built, certain components of the rulesets may be dropped at a small sacrifice to consistency and/or completeness with regard to the training data. In return, rule simplicity and a better overall performance may be gained.

The following sections discuss AQ's rule evaluation and selection process and the way the $Q(w)$ criterion can be applied at each of the above points. Section 7 presents the method that has actually been implemented in AQ18.

6.1. *Star Generation*

In the standard procedure, star generation is the process of generating a set of maximally general consistent hypotheses that cover a given positive example (a *seed*). This process involves generalizing the seed in different ways so that the obtained generalizations do not cover negative examples. This is done by applying the *extension-against* generalization rule (Michalski, 1983). The original method works sequentially, eliminating negative examples stepwise, until none are covered.

One method for implementing the $Q(w)$ measure in this process is to extend the seed against a selected group of the negative examples, and choosing the extension (generalization) with the highest $Q(w)$. This generalization is then extended against another group of negative examples, and so on, until all the negative examples are exhausted. Another method is to extend the seed sequentially against negative examples, and keep generalizations that improve $Q(w)$. This method may terminate without going to the end of the list of negative examples. The latter method is described in detail in Section 7.

6.2. *Star Completion and Rule Selection*

The $Q(w)$ measure can also be introduced at the final step of star generation when one obtained a set of rules and seeks the best rule to add to the ruleset. At this point, the candidate rules can be generalized by applying different generalization operators and selecting the best rule obtained this way. The generalization operators include dropping conditions, closing intervals, or extending references in the conditions (Michalski, 1983). The best rule from the set of generalizations is then selected using LEF.

6.3. *Rule Trimming*

Rules produced by the above star generation process are maximally general for the given coverage and consistency. It is possible, however, to change the level of generality of a rule, while preserving its coverage and consistency. AQ18 allows a user to decrease the rule generality while maintaining consistency and completeness levels by applying an operator called trim.

During the rule trimming process one can also apply a generalization operator to see if such an action would increase the $Q(w)$ value. This process can produce generalizations equivalent to those described in Section 6.2, but is more efficient because it is applied to one rule rather than a set of candidate rules. It also has a disadvantage, as it may not be able to discover generalizations that can be produced from candidate rules that were not selected during star completion.

6.4. *Rule Truncation*

Once an initial ruleset is determined, it can be truncated (or pruned) in different ways to potentially further improve it in terms of predictive accuracy or $Q(w)$ value (post-pruning). We have experimented with two post-pruning methods that involve selecting and deleting entire rules (post-specialization). One method removes low coverage rules (Zhang and Michalski, 1989). The other method removes rules that cover positive examples that can be covered by flexible matching (Kaufman and Michalski, 1999).

To introduce $Q(w)$ measure to the post-pruning process, one can apply generalization operators to the ruleset. For example, the condition dropping operator can be applied to conditions that cover many negative examples. The resulting rule will generally have lower consistency, but may have a much higher coverage. After performing such operations, truncation operators may be applied to further simplify the ruleset.

7. Implementation of $Q(w)$ in AQ18

7.1. Method

The AQ18 learning program operates in two modes (Michalski, 1999; Kaufman and Michalski, 1999). The default mode is the “noisy” mode, which relaxes the rule consistency requirement, replacing it with a $Q(w)$ formula. In the special or “no-noise” mode, AQ18 accepts only fully consistent rules, and creates a complete cover.

For initial implementation, we chose to apply the $Q(w)$ criterion as described in Sections 6.1 and 6.2: within and at the end of the star generation process. The rationale for this included ease of implementation without having to make major modifications to the existing algorithms and data structures, and the fact that this would still be relatively early in the rule selection process; thus there would remain some diversity in the set of candidate hypotheses, and it would be possible that a superior generalization of an initially inferior rule might be found.

Initially, the user was not able to modify the Q weight; instead it was automatically assumed to be 0.5 (equal emphasis on coverage and training accuracy). Thus, the user will not be bothered with having to select a particular weighting factor. In addition, this introduced the computational simplicity of not having to perform exponentiation. Since $(a^{-5} * b^{-5})$ increases monotonically with $a * b$ (assuming nonnegative a and b , of course), we replaced the strict formula for $Q(.5)$ given in Section 4 with the simpler $\text{compl}(R) * \text{consig}(R)$, thus maintaining the same ordering relationships among candidate rules.

In later implementations, we allowed the user to set a Q weight between 0 (inclusive) and 1 (exclusive), with the default remaining at 0.5. A short-cut in the code avoids the exponentiation during intermediate steps for the specific case when the weight is equal to 0.5. An example of the results of varying the Q weight on a medical dataset is shown in Section 7.2.

During star generation, AQ18 uses a beam search strategy to find the “best” generalizations of a “seed” example by a repeated application of the “extension-against” generalization rule (Michalski, 1983). In the “noisy” mode, the system determines the Q value of the generated rules after each extension-against operation; the rules with $Q(w)$ lower than that of the parent rule (the rule from which they were generated through specialization), are discarded. If the $Q(w)$ value of all rules stemming from a given parent rule is lower, the parent rule is retained instead; this operation is functionally equivalent to discarding the negative example extended against as noise.

In order to speed up the star generation, the user may specify a time-out threshold on the extension-against process. If after a given number of consecutive extensions, there has been no further improvement in rule quality, the system considers the current ruleset of sufficient quality, and terminates the extension process.

In the star termination step (i.e., after the last extension-against operation), the candidate rules are generalized additionally to determine if the resulting rules have a higher $Q(w)$ value through a hill-climbing method. Given

a rule, it tries to generalize the rule by generalizing each of its component conditions, selecting the highest-quality rule from among those generalizations, until no generalization creates further improvement.

This generalization step takes into consideration the type of the attributes in the rules, as described in (Michalski, 1983). Conditions with nominal (unordered) attributes are generalized by applying the *condition dropping* generalization operator. Conditions with linear attributes (rank, interval, cyclic, or continuous) are generalized by applying the condition dropping, the *interval extending* and the *interval closing* generalization operators. Conditions with structured attributes (hierarchically ordered) are generalized by applying the condition dropping and the *generalization tree climbing* operators. As a result of this optimization, the best rule in the resulting star is selected for output through the LEF process.

Examples of the application of these generalization rules to the base rule $[color = red \vee blue] \ \& \ [length = 2..4 \vee 10..16] \ \& \ [animal_type = dog \vee lion \vee bat]$ are presented in Table 2. In the base rule, *color* is a nominal attribute, *length* is a linear attribute, and *animal_type* is a structured attribute.

Experiments with AQ18 in this new mode exposed one unexpected difficulty. AQ18 learns rules in a “separate-and-conquer” fashion. It selects a positive example of the class it is learning, and finds the best rule it can that covers that example. If there exist examples of that class that were not covered by the selected rule, a new seed example is selected from the set of uncovered examples. A rule covering that example is added to the set of output rules, and the process repeats until all of the training examples of the positive class have been covered by some rule. Any superfluous rules (rules that do not cover any training examples that are not covered by some other rule) are then removed from the rule set, and the remaining rules are output.

Consider the following scenario. AQ learns a consistent rule that covers a seed example: $[x = 1] \ \& \ [y = 2]$. Assuming that both attributes are nominal, it attempts the generalizations by dropping conditions, thereby generating candidate rules $[x = 1]$ and $[y = 2]$. Both are found inferior to the original rule, so the originally learned rule is added to the list of rules to be output.

A subsequent seed example results in the generation of the consistent rule $[x = 1] \ \& \ [y = 5]$. Now generalizations are tested, and the rule $[x = 1]$ is found to be superior to the original rule. It is thus added to the list of rules to be output.

During the postprocessing of the list of rules, it is discovered that the rule $[x = 1] \ \& \ [y = 2]$ covers no examples uniquely (this has to be the case, since the portion of the event space covered by it is a subset of the portion of the event space covered by the rule $[x = 1]$). Hence, the rule is discarded, in favor of the rule $[x = 1]$, which was already determined to have been inferior.

To reduce the possibility of this problem occurring, we modified the Q-measuring algorithm, using a technique already present in AQ18. The options for rule preference criteria that may be implemented in a lexicographical evaluation function include two criteria for maximizing coverage. One criterion simply maximizes the number of positive events covered by a rule, while the other maximizes the coverage of positive events *that have not been covered by any previously generated and retained rule for that class*. The latter criterion, which tends toward the creation of smaller and less overlapping rulesets, has been adapted as a default, while the former one is a seldom used option.

By taking a similar course of action in Q measurement, it is possible to reduce the likelihood of a rule being supplanted by an inferior generalization. This has been implemented by changing the measurement of a rule's completeness to:

$$\text{ncmpl}(R) = p_{\text{new}} / P \quad (13)$$

where p_{new} is the number of positive events newly covered by the candidate rule.

Testing of the features described in this section indicated that the rules generated were more general than those in which every negative example was considered. Completeness was significantly higher, while consistency was slightly lower than for the analogous rules in the control data set. In addition, processing time was reduced substantially.

7.2. An Example of Results

We applied AQ18 to a medical dataset consisting of 32 attributes (all but 5 of which were of Boolean type) describing individuals' lifestyles and whether or not they had been diagnosed with various diseases. Experiments were run with three different Q weights: 0.25, 0.5, and 0.75. For the decision class *arthritis is yes*, the training set consisted of approximately 16% positive examples ($P = 1171$, $N = 6240$). The primary rules AQ18 learned in each of the three modes were as follows:

$w = 0.25$: [education \leq vocational] & [years in neighborhood > 26] &
 [rotundity \geq low] & [tuberculosis = no]:
 $p = 271$, $n = 903$, $Q = 0.111011$

$w = 0.5$: [high blood pressure = yes] & [education \leq college grad]
 $p = 355$, $n = 1332$, $Q = 0.136516$

$w = 0.75$: [education \leq college grad]
 $p = 940$, $n = 4529$, $Q = 0.303989$

As expected, increasing the Q weight allows for the reporting of rules with higher completeness and lower consistency. All three rules find a relationship between educational level and the occurrence of arthritis; the first one specializes the acceptable range of educational levels in comparison to the other two. Furthermore, one may notice that the third rule is a generalization of the second one.

8. Learning Incomplete Rulesets

The $Q(w)$ measure has been specifically developed to enable AQ-type learning methods to search for strong patterns in datasets. The previous sections focused on the relaxing the consistency requirement in determining AQ rules. This section considers the problem of relaxing the completeness requirement.

When we applied AQ18 to the medical dataset described above, the decision attributes that resulted in rulesets with any sort of strong patterns followed similar patterns in the breakdown of their rulesets. An illustrative example is one complete set of rules for determining the occurrence of arthritis. Once again, the training set consisted of 1171 positive examples (respondents who reported arthritis) and 6240 negative examples (those who did not). This was a set of examples randomly selected from the larger database, and representative of the overall distribution of positive and negative examples. When applied to this dataset, AQ18 generated a complete cover (ruleset) for the positive class (350 rules) and another for the negative class (314 rules). The

distribution of coverages of each rule in these two rulesets is shown in Tables 3a and 3b, respectively. (For the sake of brevity, some values have been combined in Table 3b into ranges).

The pattern of a dropoff in rule coverage occurred consistently in this dataset. In such cases, a complete ruleset typically contains a few strong rules and many spurious ones. Generating such a ruleset takes a significant amount of computation, which can be avoided by allowing the system to generate an incomplete ruleset.

The simplest approach to this problem is to generate only one rule for each decision class, which hopefully is the strongest. A problem with this approach is that the seed example may not be a part of the strongest rule, or may even be an error. Another problem is that this method will not detect alternative strong patterns that may exist in the data set for the given class. To overcome this problem, one may request AQ to learn n rules, where n is a fixed parameter.

The method implemented in AQ18 does not have a deterministic stopping point (a fixed n), although the user can influence the number of rules learned. The idea is to continue generating rules until sufficient performance degradation is detected. The performance of a rule is measured by the total number of positive examples covered by the rule. Two parameters define the termination condition. The first is a tolerance level, similar to that used by the LEF. If t_s is the maximum number of events covered by a stored rule thus far, the new rule must have a coverage level at list within the tolerance percentage of t_s in order to be considered a strong rule. The second parameter defines the number of consecutive weak rules that must be generated for the program to terminate. For instance, if the tolerance parameter is 20% and the iteration parameter is 3, AQ18 will continue to generate rules for a decision class until either a complete ruleset is generated, or three consecutive generated rules failed to cover at least 80% of the number of positive events covered by the rule in the ruleset with the highest coverage.

When this method was applied to the dataset described in this section, the process terminated after learning four rules for each class (as it happens, the minimum given these parameters, since it happened in this dataset that the first rule learned for each class was by far the strongest. This does not have to be the case; using this method on another decision attribute in this domain yielded seven rules for the positive class, as the strongest pattern did not appear for several iterations). In both cases, the strongest rule for each class was found. For the positive class, the third best rule from the complete ruleset and two spurious rules were also reported. For the negative class, the third, 40th, and 77th strongest of the 314 rules were returned.

It should be noted that in this mode, some of the rules returned were slight generalizations of their analogues in the complete ruleset, rather than being identical to them. This appears reasonable, as given that there are fewer rules in the final ruleset, the ones that are present need to take responsibility for more of the event space when possible.

9. Summary

This paper presented a method for integrating completeness with a newly introduced measure of consistency gain into one general measure of description (rule) quality $Q(w)$. The description consistency gain indicates the increase of classification accuracy over a random choice. The $Q(w)$ measure can be specialized to a range of criteria that weigh differently the completeness and consistency gain by varying the w parameter. The $Q(w)$

measure has been implemented in the AQ18 learning and pattern discovery system during the star generation and star termination processes.

By observing changes of the $Q(w)$ value in the process of rule generation, one can detect negative examples that may represent noise. This mechanism is particularly useful for data mining applications. By ignoring such negative examples (which increases inconsistency), the system often produces rules with much higher completeness than when it is required to satisfy the rule consistency condition. Another benefit from this mechanism is a much higher efficiency, which allows the algorithm to scale up to much more complex problems. AQ18 has also a pre-pruning mechanism for directly generating incomplete rulesets, in contrast to the post-pruning mechanism described by Bergadano et al. (1992). The pre-pruning method is based on the observation that rules selected from the first stars generated tend to be the strongest. It was also shown that in addition to completeness and consistency, a third factor, rule simplicity, can be integrated into a measure of description quality via a lexicographical evaluation functional (LEF).

Experiments have shown that the presented methodology is highly flexible and offers a powerful new tool for data mining and knowledge discovery.

References

- Baim, P.W., "The PROMISE Method for Selecting Most Relevant Attributes for Inductive Learning Systems," Report No. UIUCDCS-F-82-898, Department of Computer Science, University of Illinois, Urbana, 1982.
- Bergadano, F., Matwin, S., Michalski R.S. and Zhang, J., "Learning Two-tiered Descriptions of Flexible Concepts: The POSEIDON System," *Machine Learning* 8, 1992, pp. 5-43.
- Bruha, I., "Quality of Decision Rules: Definitions and Classification Schemes for Multiple Rules," In Nakhaeizadeh, G. and Taylor, C.C. (eds.), *Machine Learning and Statistics, The Interface*, New York: John Wiley & Sons, Inc., 1997, pp. 107-131.
- Clark, P. and Boswell, R., "Rule Induction with CN2: Some Recent Improvements," in Kodratoff, Y. (ed.), *Proceedings of the Fifth European Working Session on Learning (EWSL-91)*, Berlin: Springer-Verlag, 1991, pp. 151-163.
- Clark, P. and Niblett, T., "The CN2 Induction Algorithm," *Machine Learning* 3, 1989, pp. 261-283.
- Cohen, W., "Fast Effective Rule Induction," *Proceedings of the Twelfth International Conference on Machine Learning*, Lake Tahoe, CA, 1995.
- Fayyad, U.M., Piatetsky-Shapiro, G., Smyth, P. and Uthurusamy, R. (eds.), *Advances in Knowledge Discovery and Data Mining*, Menlo Park, CA: AAAI Press, 1996.
- Fürnkranz, J. and Widmer, G., "Incremental Reduced Error Pruning," *Proceedings of the Eleventh International Conference on Machine Learning*, New Brunswick, NJ, 1994.
- Kaufman, K.A., "INLEN: A Methodology and Integrated System for Knowledge Discovery in Databases," Ph.D. dissertation, *Reports of the Machine Learning and Inference Laboratory*, MLI 97-15, George Mason University, Fairfax, VA, 1997.
- Kaufman, K.A. and Michalski, R.S., "STAR: An Environment for Natural Induction and Learning," *Reports of the Machine Learning and Inference Laboratory*, George Mason University, Fairfax, VA, 1999 (to appear).
- Michalski, R.S., "A Theory and Methodology of Inductive Learning," In Michalski, R.S. Carbonell, J.G. and Mitchell, T.M. (eds.), *Machine Learning: An Artificial Intelligence Approach*, Palo Alto: Tioga Publishing, 1983, pp. 83-129.
- Michalski, R.S., "NATURAL INDUCTION: A Theory, Methodology and Applications to Knowledge Mining and Pattern Discovery," *Reports of the Machine Learning and Inference Laboratory*, George Mason University, 1999 (to appear).
- Piatetsky-Shapiro, G., "Discovery, Analysis, and Presentation of Strong Rules," in Piatetsky-Shapiro, G. and Frawley, W. (eds.), *Knowledge Discovery in Databases*, Menlo Park, CA: AAAI Press, 1991, pp. 229-248.
- Quinlan, J.R., "Induction of Decision Trees," *Machine Learning* 1, pp. 81-106, 1986.
- Zhang, J. and Michalski, R.S., "Rule Optimization via SG-TRUNC Method," *Proceedings of the Fourth European Working Session on Learning*, Montpellier, 1989, pp. 251-262.

Table 1. A comparison of rule evaluation criteria.

Data Set	Pos	Neg	Info		Gain		PROMISE		CN2		RIPPER		Q(0)		Q(.25)		Q(.5)		Q(.75)		Q(1)	
			V	R	V	R	V	R	V	R	V	R	V	R	V	R	V	R	V	R	V	R
A	200	pos	50	5	.10	7	.24	7	.44	4	.05	7	.89	4	.65	7	.47	7	.34	7	.25	6
			50	0	.12	6	.25	6	0	1	.05	6	1	1	.71	6	.5	6	.35	6	.25	6
		200	5	.69	1	.99	1	.17	2	.20	1	.97	2	.98	1	.99	1	.99	1	1	1	
	800	pos	150	10	.39	2	.74	2	.34	3	.14	2	.92	3	.88	2	.83	2	.79	2	.75	2
			150	30	.33	3	.71	3	.65	6	.12	3	.79	6	.78	3	.77	3	.76	3	.75	2
		neg	100	15	.21	5	.48	5	.56	5	.09	5	.84	5	.74	4	.65	4	.57	5	.5	5
			120	25	.24	4	.57	4	.66	7	.10	4	.78	7	.73	5	.69	5	.64	4	.6	4
B	500	pos	50	5	.03	7	.09	7	.44	3	.05	7	.82	3	.48	7	.29	7	.17	7	.1	7
			250	25	.21	6	.45	5	.44	3	.23	5	.82	3	.72	5	.64	5	.57	5	.5	5
		500	50	.76	1	.9	1	.44	3	.45	1	.82	3	.86	1	.91	1	.95	1	1	1	
	500	pos	500	150	.49	2	.7	3	.78	7	.35	3	.54	7	.63	6	.73	4	.86	2	1	1
			200	5	.21	5	.39	6	.17	1	.20	6	.95	1	.77	4	.62	6	.5	6	.4	6
		neg	400	35	.44	3	.73	2	.40	2	.37	2	.84	2	.83	2	.82	2	.81	3	.8	3
			400	55	.38	4	.69	4	.53	6	.35	4	.76	6	.77	3	.78	3	.79	4	.8	3
C	800	pos	50	5	.004	7	0	-	.44	3	.05	7	.55	3	.32	6	.18	6	.11	6	.06	7
			250	25	.02	5	0	-	.44	3	.23	5	.55	3	.47	4	.41	5	.36	4	.31	5
		500	50	.07	1	0	-	.44	3	.45	1	.55	3	.56	3	.58	1	.60	1	.63	1	
	200	pos	500	150	.01	6	0	-	.78	7	.35	3	<0	7	<0	7	<0	7	<0	7	.63	1
			200	5	.05	3	0	-	.17	1	.20	6	.88	1	.64	1	.47	3	.34	5	.25	6
		neg	400	35	.05	2	0	-	.40	2	.37	2	.6	2	.57	2	.55	2	.52	2	.5	3
			400	55	.02	4	0	-	.53	6	.35	4	.4	6	.42	5	.44	4	.47	3	.5	3

Columns labeled *V* indicate raw value.

Columns labeled *R* indicate rank assigned by the given evaluation method in the given dataset.

Table 2. Effects of different Q-optimizing generalization operators on the base rule.

Generalization Type	Resulting Rule
Removing nominal condition	$[length = 2..4 \vee 10..16] \ \& \ [animal_type = dog \vee lion \vee bat]$
Removing linear condition	$[color = red \vee blue] \ \& \ [animal_type = dog \vee lion \vee bat]$
Extending linear interval	$[color = red \vee blue] \ \& \ [length = 2..4 \vee 8..16] \ \& \ [animal_type = dog \vee lion \vee bat]$
Closing linear interval	$[color = red \vee blue] \ \& \ [length = 2..16] \ \& \ [animal_type = dog \vee lion \vee bat]$
Removing structured condition	$[color = red \vee blue] \ \& \ [length = 2..4 \vee 10..16]$
Generalizing structured condition	$[color = red \vee blue] \ \& \ [length = 2..4 \vee 10..16] \ \& \ [animal_type = mammal]$

Base rule: $[color = red \vee blue] \ \& \ [length = 2..4 \vee 10..16] \ \& \ [animal_type = dog \vee lion \vee bat]$

Table 3. Distribution of coverage levels in the Arthritis ruleset.

a. For the positive class

b. For the negative class

Rule Coverage	Number of Rules
325	1
126	1
55	1
45	1
19	3
17	1
13	1
12	2
10	3
9	1
8	2
7	4
6	6
5	9
4	11
3	37
2	67
1	199
Totals: 1171 Positive Examples	350 Rules

Rule Coverage	Number of Rules
3092	1
300-500	5
250-299	4
200-249	6
150-199	10
100-149	7
75-99	6
50-74	14
25-49	36
15-24	38
10-14	25
8-9	14
6-7	14
5	14
4	18
3	20
2	28
1	54
6240 Negative Examples	314 Rules