

Two-Way Induction

Pedro Domingos
Department of Information and Computer Science
University of California, Irvine
Irvine, California 92717, U.S.A.
pedrod@ics.uci.edu

Abstract

General-to-specific learners like ID3 and CN2 perform well when the target concept descriptions are general, but often have difficulties when they are specific or mixed. This problem can be alleviated by combining them with a specific-to-general learning component, resulting in a two-way induction system. In this paper one design for such a component is proposed, as well as two methods for combining the two components. Experiments on artificial domains show the combined learner to match or outperform “pure” versions of C4.5 and CN2 across the entire generality spectrum, with the advantage increasing for greater concept specificity. Experiments on 24 real-world domains from the UCI repository confirm the utility of two-way induction: the combined learner achieves higher accuracy than C4.5 in 17 domains (at the 5% significance level in 12), and similar results are obtained with CN2. Closer observation of the system’s behavior leads to a better understanding of its ability to correct overly-general rules with specific ones, and shows that there is still room for improvement.

1 Introduction and Motivation

Most widely-used decision-tree and rule learners perform general-to-specific induction: they start with an empty concept description, and gradually add restrictions to it until there is not enough evidence to continue, or perfect discrimination is achieved. ID3/C4.5 [14] and CN2 [4] are examples of systems that function in this way. They have some significant advantages: because only statistical measures are used to evaluate induction steps, as opposed to individual examples being considered, good noise immunity can be achieved. Also, when the induction process stops early, as it often does, learning can be fast, and the resulting descriptions concise.

However, as in any search process, finding the target tends to be harder when it lies farther from the search’s starting point, in this case the null concept description. The opportunities for the search to take the wrong path accumulate, and the probability of an incorrect end result increases. As a result, if all or part of the concept description is fairly specific, i.e., if it includes a large fraction of the attributes used to describe examples, systems like C4.5 and CN2 may be unable to find it; they are prone to either stopping the search too soon, leading to an overly general description, or to inducing erroneous rules/branches that may or may not be pruned. This performance degradation as concept specificity increases is indeed verified in artificial domains in a later section of this paper.

Specific-to-general learners, on the other hand, should find it easier to induce fairly specific descriptions. They have other disadvantages, however: search has to start from specific examples, making it more sensitive to noise, and early decisions, based on little evidence, can be wrong and difficult to correct later on. As a result, even though some specific-to-general learners exist (e.g., EACH [16]), they have not come into widespread use.

More generally, when the target description is a mix of general rules and more specific “exception” areas—a pattern which appears to occur frequently in real applications—we can expect both approaches to have problems: general-to-specific learners may not recognize the exception areas, and specific-to-general ones may induce only imperfect, corrupted versions of the general rules. A natural solution would then be to combine the two search directions in a single system, allowing both the induction of rules starting from the null description and from specific examples, and employing some conflict resolution strategy in parts of the instance space where rules’ predictions disagree. This paper describes and evaluates one realization of this idea.

C4.5 and CN2 were used in turn as the general-to-specific learner. The next section describes the specific-to-general component used, and the following one describes two methods for combining the two components. The results of experiments on natural and artificial domains are then reported and interpreted. Finally related work is discussed, and some conclusions and directions for future work are presented.

2 Specific-to-General Induction

Specific-to-general induction is performed by the RISE algorithm, of which the learning and classification procedures will be considered in turn. More details can be found in [7, 6].

2.1 Representation and Search

Each example is a vector of attribute-value pairs, together with a specification of the class to which it belongs; attributes can be either nominal (symbolic) or numeric. Each rule consists of a conjunction of antecedents and a predicted class. Each antecedent is a condition on a single attribute, and there is at most one antecedent per attribute. Conditions on nominal attributes are equality tests of the form $a_i = v_j$, where a_i is the attribute and v_j is one of its legal values. Conditions on numeric attributes take the form of allowable intervals for the attributes, i.e., $a_i \in [v_{j1}, v_{j2}]$, where v_{j1} and v_{j2} are two legal values for a_i . Exemplars (i.e., examples used as prototypes for classification) are viewed as maximally specific rules, with conditions on all attributes and degenerate (point) intervals for numeric attributes. A rule is said to *cover* an example if the example satisfies all of the rule's conditions; a rule is said to *win* an example if it is the nearest rule to the example according to the distance metric that will be described below.

The RISE algorithm is summarized in Table 1. RISE searches for "good" rules in a specific-to-general fashion, starting with a rule set that is the training set of examples itself. RISE looks at each rule in turn, finds the nearest example of the same class that it does not already cover (i.e., that is at a distance greater than zero from it), and attempts to minimally generalize the rule to cover it. The generalization procedure is outlined in Table 2. If the change's effect on global accuracy is positive, it is retained; otherwise it is discarded. Generalizations are also accepted if they appear to have no effect on accuracy; this reflects a simplicity bias. This procedure is repeated until, for each rule, attempted generalization fails.

Table 1: The RISE algorithm.

Input: ES is the training set.

Procedure RISE (ES)

Let RS be ES .

Compute $Acc(RS)$.

Repeat

For each rule R in RS ,

Find the nearest example E to R not already covered by it (and of the same class).

Let $R' = \text{Most_Specific_Generalization}(R, E)$.

Let $RS' = RS$ with R replaced by R' .

If $Acc(RS') \geq Acc(RS)$

Then Replace RS by RS' ,

If R' is identical to another rule in RS ,

Then delete R' from RS .

Until no increase in $Acc(RS)$ is obtained.

Return RS .

A potential difficulty is that measuring the accuracy of a rule set on the training set requires matching all rules with all training examples, and this would entail a high computational cost if it was repeatedly done as outlined. Fortunately, at each step only the *change* in accuracy needs to be computed. Each example memorizes the distance to its nearest rule and its assigned class. When a rule is generalized, all that is necessary is then to match that single rule against all examples, and check if it wins any that it did not before, and what its effect on these is. Previously

Table 2: Generalization of a rule to cover an example.

Inputs: $R = (a_1, a_2, \dots, a_A, c_R)$ is a rule,
 $E = (e_1, e_2, \dots, e_A, c_E)$ is an example.
 a_i is either True, $x_i = r_i$, or $r_{i,lower} \leq x_i \leq r_{i,upper}$.

Function Most_Specific_Generalization (R, E)

For each attribute i ,

If $a_i = \text{True}$ then Do nothing.

Else if i is symbolic and $e_i \neq r_i$ then $a_i = \text{True}$.

Else if $e_i > r_{i,upper}$ then $r_{i,upper} = e_i$.

Else if $e_i < r_{i,lower}$ then $r_{i,lower} = e_i$.

misclassified examples that are now correctly classified add to the accuracy, and previously correctly classified examples that are now misclassified subtract from it. If the former are more numerous than the latter, the change in accuracy is positive, and the generalization is accepted. With this optimization, RISE’s worst-case time complexity has been shown to be quadratic in the number of examples and the number of attributes, which is comparable to that of general-to-specific rule induction algorithms [7].

2.2 Classification

At performance time, classification of each test example is performed by finding the nearest rule to it, and assigning the example to the rule’s class. The distance measure used is a combination of Euclidean distance for numeric attributes, and a simplified version of Stanfill and Waltz’s value difference metric for symbolic attributes [18].

Let $E = (e_1, e_2, \dots, e_A, c_E)$ be an example with value e_i for the i th attribute and class c_E . Let $R = (a_1, a_2, \dots, a_A, c_R)$ be a rule with class c_R and condition a_i on the i th attribute, where $a_i = \text{True}$ if there is no condition on i , otherwise a_i is $x_i = r_i$ if i is symbolic and a_i is $r_{i,lower} \leq x_i \leq r_{i,upper}$ if i is numeric. The distance $\Delta(R, E)$ between R and E is then defined as:

$$\Delta(R, E) = \sum_{i=1}^A \delta^2(i) \quad (1)$$

where the component distance $\delta(i)$ for the i th attribute is:

$$\delta(i) = \begin{cases} 0 & \text{if } a_i = \text{True} \\ SVDM(r_i, e_i) & \text{if } i \text{ is symbolic} \wedge a_i \neq \text{True} \\ \delta_{num}(i) & \text{if } i \text{ is numeric} \wedge a_i \neq \text{True} \end{cases} \quad (2)$$

$SVDM(r_i, e_i)$ is the simplified value difference metric, defined as:

$$SVDM(x_i, x_j) = \sum_{h=1}^C |P(c_h|x_i) - P(c_h|x_j)| \quad (3)$$

where x_i and x_j are any legal values of the attribute, C is the number of classes, c_h is the h th class, and $P(c_h|x_i)$ denotes the probability of c_h conditioned on x_i . The essential idea behind VDM-type metrics is that two values should be considered similar if they make similar class predictions, and dissimilar if their

predictions diverge. This has been found to give good results in several domains [5]. Notice that, in particular, $SVDM(x_i, x_j)$ is always 0 if $i = j$.

The component distance for numeric attributes is defined as:

$$\delta_{num}(i) = \begin{cases} 0 & \text{if } r_{i,lower} \leq e_i \leq r_{i,upper} \\ \frac{e_i - r_{i,upper}}{x_{max} - x_{min}} & \text{if } e_i > r_{i,upper} \\ \frac{r_{i,lower} - e_i}{x_{max} - x_{min}} & \text{if } e_i < r_{i,lower} \end{cases} \quad (4)$$

x_{max} and x_{min} being respectively the maximum and minimum observed values for the attribute.

The distance from a missing numeric value to any other is defined as 0. If a symbolic attribute’s value is missing, it is assigned the special value “?”. This is treated as a legitimate symbolic value, and its distance to all other values of the attribute is computed and used. When coupled with SVDM, this is a sensible policy: a missing value is taken to be roughly equivalent to a given possible value if it behaves similarly to it, and inversely if it doesn’t.

When two or more rules are equally close to a test example, the rule that was most accurate on the training set wins. So as to not unduly favor more specific rules, the Laplace-corrected accuracy is used [12]:

$$LAcc(R) = \frac{N_{corr}(R) + 1}{N_{won}(R) + C} \quad (5)$$

where R is any rule, C is the number of classes, $N_{won}(R)$ is the total number of examples won by R , $N_{corr}(R)$ is the number of examples among those that R correctly classifies, and C is the number of classes. The effect of the Laplace correction is to make the estimate of a rule’s accuracy converge to the “random guess” value of $1/C$ as the number of examples won by the rule decreases. Thus rules with high apparent accuracy are favored only if they also have high statistical support, i.e., if that apparent accuracy is not simply the result of a small sample.

3 Two-Way Induction

Applying the RISE algorithm to a training set produces a set of rules which will henceforth be called the S rules (for “specific”). Applying the general-to-specific learner produces a set of G rules. If CN2 is used, the rules it outputs are used directly. If C4.5 is used, C4.5RULES is first applied to convert the decision tree it produces into a set of rules [14].

Two algorithms for combining the S and G rules were designed, TWI-1 and TWI-2. In TWI-1, the sets of S and G rules are merged to form a single set, deleting any duplicates. The Laplace accuracy of each rule on the examples that it *covers* is then measured, and the classification procedure is applied to each training example in turn; each rule memorizes the number of examples it won, and how many of them it classified correctly. At the end, the Laplace accuracy of each rule on the examples that it *won* is computed, and this is the value retained. The earlier estimates were necessary to break ties during the classification cycle that produced the final ones. At classification time, the procedure described in the previous section is applied without any distinction between S and G rules. The nearest rule to the test example wins; if two or more rules are equally near, the one with the highest Laplace accuracy prevails.

In the TWI-2 algorithm, the sets of S and G rules are kept separate, and the Laplace accuracy of each rule on the training examples that it covers is measured. At classification time, the two sets of rules are first applied separately. A winner among the S rules is found by the procedure previously outlined. To select the G winner, the G rules are matched against the test example. If only one G rule covers it, that rule is chosen as the G winner. If more than one rule covers the example, the one with the highest Laplace accuracy wins. If no G rule covers the example, the default rule is chosen as the G winner. Of the two finalists (S and G), the one with the highest accuracy then wins and classifies the example.

4 Empirical Study

An empirical study was conducted to evaluate the contribution of the specific-to-general learner, and the performance of the two-way induction system relative to its individual components. A recent version of CN2 was used [3]. The default settings for C4.5, C4.5RULES and CN2 were kept. Note that all results for C4.5RULES (abbreviated C4.5R) and CN2 below are those obtained using their own classification procedures, not the one above for G rules. Experiments in artificial and real-world domains are described in turn.

4.1 Artificial Domains

The goal of the study in artificial domains was to verify in ideal, controlled conditions if the hypotheses

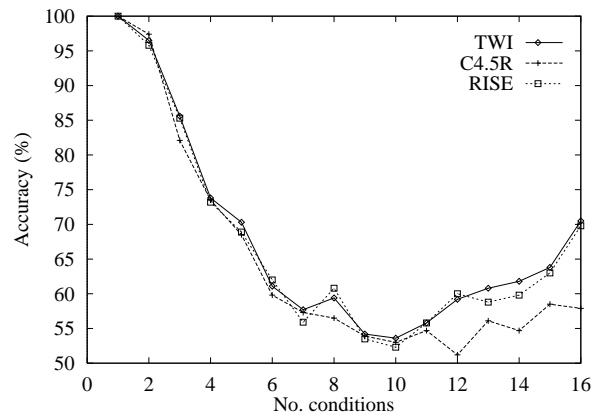


Figure 1: Accuracy as a function of concept specificity.

put forth in the introduction to this paper hold. Following a philosophy justified elsewhere [17, 1], classes of domains rather than individual concept definitions were considered. The study was carried out using C4.5R, RISE and TWI-2.

The independent variable of interest is the specificity of the target concept description. A good operational measure of it is the average length of the rules comprising the description: rules with more conditions imply a more specific concept. The dependent variables are the out-of-sample accuracies of the three algorithms. Concepts defined as Boolean functions in disjunctive normal form were used as targets. The datasets were composed of 100 examples described by 16 attributes. The average number of literals C in each disjunct comprising the concept was varied from 1 to 16. The number of disjuncts was set to $\text{Min}\{2^{C-1}, 25\}$. This attempts to keep the fraction of the instance space covered by the concept roughly constant, up to the point where it would require more rules than could possibly be learned. Equal numbers of positive and negative examples were included in the dataset, and positive examples were divided evenly among disjuncts. In each run a different target concept was used, generating the disjuncts at random, with length given by a binomial distribution with mean C and variance $C(1 - \frac{C}{16})$; this is obtained by including each feature in the disjunct with probability $\frac{C}{16}$. Twenty runs were conducted, with two-thirds of the data used for training and the remainder for testing.

The results are shown graphically in Fig. 1. The most salient aspect is the large difference in difficulty between short and long rules for all learners. Concepts with very few (approx. three or less) conditions per

Table 3: Experimental results when G is C4.5.

Domain	TWI-2	TWI-1	C4.5	RISE	Default
Audiology	78.9±4.3	66.3±6.6 ¹	70.6±5.7 ¹	77.3±4.9 ¹	21.3±2.6 ¹
Breast cancer	69.4±4.5	67.5±5.5 ⁴	67.8±6.6 ⁵	68.2±4.2 ³	67.6±7.6 ⁶
Credit	83.5±1.7	84.2±2.9 ⁶	84.8±2.5 ³	83.1±2.0 ⁴	57.4±3.8 ¹
Echocardiogram	68.1±4.9	67.3±8.2 ⁶	65.9±7.6 ⁶	67.4±4.9 ⁵	67.8±6.6 ⁶
Glass	70.2±6.5	64.1±10.1 ²	64.9±9.2 ¹	69.7±6.1 ⁶	31.7±5.5 ¹
Heart disease	79.8±3.7	76.6±3.1 ¹	77.8±4.3 ⁴	79.6±3.9 ⁶	55.0±3.4 ¹
Hepatitis	78.1±5.4	78.8±3.7 ⁶	78.6±5.3 ⁶	76.9±5.3 ³	78.1±3.1 ⁶
Horse colic	83.8±3.5	81.2±6.3 ⁵	81.2±4.4 ²	81.9±3.2 ¹	63.6±3.9 ¹
Iris	93.0±2.7	93.6±2.7 ⁶	93.2±2.5 ⁶	92.9±2.8 ⁶	26.5±5.2 ¹
LED	69.5±4.0	69.1±3.8 ⁶	69.0±3.8 ⁶	67.9±3.6 ³	9.9±3.0 ¹
Labor	83.7±9.9	81.6±10.7 ⁶	86.3±8.6 ⁶	89.2±10.6 ³	65.0±9.5 ¹
Liver disease	64.0±5.6	65.8±4.6 ⁶	64.4±3.9 ⁶	63.4±5.4 ⁴	58.1±3.4 ¹
Lung cancer	48.6±15.1	40.5±14.0 ⁴	40.5±14.0 ⁴	50.5±15.2 ⁵	26.8±12.3 ¹
Lymphography	78.4±6.2	76.9±3.9 ⁶	75.6±4.9 ⁴	80.2±6.8 ¹	57.3±5.4 ¹
Pima diabetes	70.8±3.0	73.0±3.1 ²	74.3±3.0 ¹	70.4±3.0 ⁴	66.0±2.3 ¹
Post-operative	67.8±6.0	71.0±5.2 ²	68.2±6.9 ⁶	62.3±9.1 ¹	71.2±5.2 ²
Primary tumor	41.4±4.6	35.4±5.4 ¹	37.5±5.7 ¹	39.8±5.3 ¹	24.6±3.2 ¹
Promoters	87.7±5.8	80.6±10.1 ¹	80.4±10.0 ¹	87.7±5.5 ⁶	43.1±4.2 ¹
Solar flare	71.9±3.1	71.2±4.1 ⁶	71.1±4.1 ⁶	70.8±2.7 ²	25.2±4.4 ¹
Sonar	74.7±12.1	64.3±9.4 ¹	65.4±7.1 ¹	76.0±11.4 ³	50.8±7.6 ¹
Soybean	83.1±6.8	75.1±5.8 ¹	78.9±5.9 ¹	84.8±6.5 ¹	9.1±2.1 ¹
Voting records	95.9±1.6	95.7±1.7 ⁶	95.8±1.3 ⁶	95.5±1.5 ⁵	60.5±3.1 ¹
Wine	96.3±2.3	91.3±5.6 ¹	91.8±5.6 ¹	96.9±1.8 ⁵	36.4±9.9 ¹
Zoology	93.5±3.8	90.0±5.2 ¹	89.6±4.7 ¹	93.2±3.7 ⁶	39.4±6.4 ¹

rule are so simple that both types of learner are able to learn them easily; the expected advantage of the general-to-specific approach is not realized. In separate experiments, corrupting the data with 10% and 20% noise degraded the performance of all algorithms equally, again giving no advantage to C4.5R. At the other end, however, specific-to-general learning has a clear advantage for concepts with 12 or more conditions per rule; all differences here are significant at the 5% level using a one-tailed paired t test. (Error bars are not shown to avoid further cluttering the graph.) The combined approach is able to effectively leverage both components, even when one is substantially less accurate than the other; TWI performs better than the best of C4.5 and RISE in 10 of the 16 data points, including better than RISE at the 5% error level for $C = 13$ and $C = 14$.

The slight upward trend in C4.5R’s curve for $C > 10$ was investigated by repeating the experiments with 32 attributes, 400 examples, a maximum of 50 rules and $C = 1, \dots, 32$. C4.5R’s lag increases, but the upward trend is maintained; on inspection of the rules C4.5R produces, this is revealed to be due to the

fact that, as the concept rules become more and more specific, it becomes possible to induce short rules for its negation. The hardest concepts, for which both the concept and its negation have necessarily long rules, are for intermediate values of C .

4.2 Natural Domains

Whether the expected advantages of adding a specific-to-general component to a general-to-specific learner are realized in real-world applications was investigated by carrying out experiments on 24 domains from the UCI repository [11]. Again, 20 runs were conducted for each dataset, with two-thirds of the data used for training. The default classifier (assigning examples to the most frequent class) was included as a baseline. The results are shown in Table 3 for C4.5 used as the G component, and in Table 4 for CN2. The average accuracy and sample standard deviation are tabulated for each algorithm in each domain. The superscripts denote significance levels for the difference in accuracy between TWI-2 and the corresponding algorithm, using a one-tailed paired t test: 1 denotes 0.5%, 2 is 1%, 3 is 2.5%, 4 is 5%, 5 is 10%, and 6 is

Table 4: Experimental results when G is CN2.

Domain	TWI-2	TWI-1	CN2	RISE	Default
Audiology	77.3±4.4	63.9±7.4 ¹	71.0±5.1 ¹	77.3±4.9 ⁶	21.3±2.6 ¹
Breast cancer	68.4±6.3	67.9±6.4 ⁶	67.9±7.1 ⁶	68.2±4.2 ⁶	67.6±7.6 ⁶
Credit	84.4±2.2	82.3±2.2 ¹	82.0±2.2 ¹	83.1±2.0 ¹	57.4±3.8 ¹
Echocardiogram	68.0±5.3	66.0±5.9 ³	68.2±7.2 ⁶	67.4±4.9 ⁶	67.8±6.6 ⁶
Glass	69.9±6.8	58.9±6.3 ¹	63.8±5.5 ¹	69.7±6.1 ⁶	31.7±5.5 ¹
Heart disease	81.2±4.3	79.5±2.9 ³	79.7±2.9 ⁴	79.6±3.9 ²	55.0±3.4 ¹
Hepatitis	79.6±5.3	80.5±4.5 ⁶	80.3±4.2 ⁶	76.9±5.3 ¹	78.1±3.1 ⁵
Horse colic	83.1±3.6	83.0±3.6 ⁶	82.5±4.2 ⁶	81.9±3.2 ¹	63.6±3.9 ¹
Iris	93.4±2.7	93.2±2.8 ⁶	93.3±3.6 ⁶	92.9±2.8 ⁴	26.5±5.2 ¹
LED	69.5±4.0	68.9±4.5 ⁶	69.5±3.8 ⁶	67.9±3.6 ¹	9.9±3.0 ¹
Labor	87.4±10.6	82.9±8.9 ³	82.1±6.9 ⁴	89.2±10.6 ⁶	65.0±9.5 ¹
Liver disease	65.9±5.2	66.4±3.8 ⁶	65.0±3.8 ⁶	63.4±5.4 ¹	58.1±3.4 ¹
Lung cancer	42.3±14.5	39.1±14.8 ⁶	38.6±13.5 ⁶	50.5±15.2 ¹	26.8±12.3 ¹
Lymphography	80.9±6.3	79.4±5.0 ⁵	78.8±4.9 ⁴	80.2±6.8 ⁶	57.3±5.4 ¹
Pima diabetes	72.3±2.7	73.7±2.9 ²	73.8±2.7 ¹	70.4±3.0 ¹	66.0±2.3 ¹
Post-operative	68.3±4.9	65.0±6.9 ³	60.8±8.2 ¹	62.7±9.2 ¹	71.2±5.2 ¹
Primary tumor	41.4±5.2	38.9±5.1 ²	39.8±5.2 ⁴	39.8±5.3 ¹	24.6±3.2 ¹
Promoters	85.6±6.2	78.1±9.1 ¹	75.9±8.8 ¹	87.7±5.5 ⁵	43.1±4.2 ¹
Solar flare	71.4±2.6	69.6±3.6 ²	70.4±3.0 ⁵	70.8±2.6 ⁶	25.2±4.4 ¹
Sonar	73.8±10.6	63.2±9.0 ¹	66.2±7.5 ¹	76.0±11.4 ³	50.8±7.6 ¹
Soybean	83.5±6.7	77.8±6.7 ¹	77.4±7.2 ¹	84.8±6.5 ¹	9.1±2.1 ¹
Voting records	95.0±1.8	95.5±1.5 ⁴	95.8±1.6 ¹	95.5±1.5 ⁴	60.5±3.1 ¹
Wine	95.3±2.6	91.0±4.5 ¹	90.8±4.7 ¹	96.9±1.8 ¹	36.4±9.9 ¹
Zoology	93.7±3.3	90.9±5.1 ¹	90.6±5.0 ¹	93.2±3.7 ⁵	39.4±6.4 ¹

above 10%. For example, in the horse colic domain in Table 3, TWI-2 obtains an accuracy of 83.8% with a standard deviation of 3.5%, is better than C4.5 at the 1% level, and better than RISE at the 0.5% level.

Table 5 summarizes these results. The first line shows the number of domains in which TWI-2 achieved higher accuracy than the corresponding system, vs. the number in which the reverse happened. The second line considers only those domains in which the observed difference is significant at the 5% level, and shows that most of the previous “wins” were indeed significant. For example, when G was C4.5, TWI-2 did better than C4.5 in 17 domains overall, and worse in 7; with 5% significance it did better in 12 and worse in 2. All other comparisons are similarly favorable to TWI-2. The third line shows the results of applying a sign test to the values of line one. This consists of considering the number of wins as a binomial variable, and asking how unlikely the value obtained is if TWI-2 is assumed to be no more accurate than the corresponding algorithm. For example, 17 wins in 24 trials has a probability of occurrence of only 3%. This test shows that all the numbers of wins obtained

are significant at the 5% level, resulting in high confidence that TWI-2 is a more accurate learner than either component alone, if the set of domains used is assumed to be representative of real-world tasks.

4.3 Discussion

The utility of two-way induction is clearly shown by the results above, but it is important to understand why this advantage is observed. Inspection of the S and G rules shows that the S rules are indeed substantially more specific than the G ones. G rules typically contain a small number of conditions (approx. 1 to 5), whereas S rules often contain conditions on half of all the attributes. Further, tracing of TWI-2 reveals that by far the majority of S wins occur when the test example is also covered by a G rule, but G rules are still correct on most of the examples they cover. S rules thus encapsulate small exception areas to the broad generalizations represented by G rules, as intended.

TWI-2 performed better than TWI-1, and this was to be expected. TWI-1’s naive method of combining S and G rules (merge all) has several disadvantages.

Table 5: Summary of TWI-2’s results vs. other algorithms.

Measure	G = C4.5			G = CN2		
	TWI-1	C4.5	RISE	TWI-1	CN2	RISE
No. wins	18-6	17-7	17-6	20-4	19-4	16-7
No. signif. wins	11-2	12-2	11-4	14-2	13-2	10-5
Sign test	1.0	3.0	2.0	0.1	0.2	5.0

The G rules were not designed to be applied in a best-match manner. Being very general, they easily win examples outside them over more appropriate S rules. The S rules, conversely, were designed for best-match classification, but in the large sections of the instance space covered by G rules, they can only win examples that they actually cover.

TWI-2’s accuracy is bounded from below by the fraction of cases on which the S and G components agree and are correct, and from above by 100% minus the fraction of cases on which they agree and are both incorrect. Further, when there are more than two classes, S and G components may differ and both be wrong. Extracting the cases where TWI-2’s decision makes a difference, we see that it makes the correct decision approximately two-thirds of the time, on average. This is well above chance, but still leaves substantial room for improvement, at least in theory.

The study above was mainly concerned with accuracy. Two other variables of interest are speed and comprehensibility of the results. In practice, RISE is slower than C4.5 and CN2, but this is only noticeable in larger datasets (more than 500 examples). Also, because the S rules are longer, they are harder for a human to understand; but since in TWI-2 the distinction between S and G rules is maintained, the G rules can be seen as an approximate and accessible model of the domain, while the S rules represent second-order refinements and exceptions that should be taken into account to achieve higher accuracy.

5 Related Work

Mitchell’s [10] version space approach is an early example of two-way induction. TWI is a heuristic algorithm, in contrast to Mitchell’s exhaustive candidate elimination procedure, but has the advantage that its worst-case space and time complexities are respectively linear and low-order polynomial, instead of exponential. Also, unlike version spaces, TWI is able to deal with disjunctive concepts, noise, and missing and continuous attributes.

The research described in this paper falls in the general area of empirical multi-strategy learning, which attempts to produce more accurate learners by combining two or more individual algorithms [9]. The MCS system combines decision trees with the nearest-neighbor algorithm and linear machines [2]. FCLS combines rules with specific examples in a best-match framework [20]. Quinlan has combined decision trees and other methods with nearest-neighbor in regression tasks [15]. The post-pruning step used in many rule and decision-tree learners can be seen as a form of specific-to-general induction, and Quinlan [13] and others have shown it to be effective. Specific-to-general induction is performed by systems like NGE [16] and KBNGE [19], but neither combines it with general-to-specific learning. Golding and Rosenbloom’s Anapron system for name pronunciation [8] performs search in neither direction, but is similar in spirit to TWI: a set of expert-supplied rules functions as the G component, and a set of implicit rules (the “arules”) contained in a case base functions as the S component.

6 Conclusions and Future Work

This paper has shown that the accuracy of general-to-specific learners like ID3/C4.5 and CN2 can be substantially improved by adding a specific-to-general component to them. As the observations made indicate, there is still room for improvement, and work on more sophisticated methods of combining the two components is underway. Inputting unpruned rules or trees and combining the post-pruning stage with the specific-to-general induction process is also a potentially productive avenue, and preliminary work has been done. Another direction for future research is to use an expert-supplied domain theory as the G component, leading to a combined analytical-empirical learner.

Acknowledgments

This work was partly supported by JNICT/Programa Ciência and Fulbright scholarships. The author is grateful to Dennis Kibler and Mike Pazzani, to the creators of the C4.5 and CN2 systems, and to all the people who provided the datasets used in the empirical study. Please see the documentation in the UCI Repository for detailed information.

References

- [1] D. W. Aha. Generalizing from case studies: A case study. In *Proceedings of the Ninth International Workshop on Machine Learning*, pages 1–10, Aberdeen, Scotland, 1992. Morgan Kaufmann.
- [2] C. E. Brodley. Addressing the selective superiority problem: Automatic algorithm/model class selection. In *Proceedings of the Tenth International Conference on Machine Learning*, pages 17–24, Amherst, MA, 1993. Morgan Kaufmann.
- [3] P. Clark and R. Boswell. Rule induction with CN2: Some recent improvements. In *Proceedings of the Sixth European Working Session on Learning*, pages 151–163, Porto, Portugal, 1991. Springer-Verlag.
- [4] P. Clark and T. Niblett. The CN2 induction algorithm. *Machine Learning*, 3:261–283, 1989.
- [5] S. Cost and S. Salzberg. A weighted nearest neighbor algorithm for learning with symbolic features. *Machine Learning*, 10:57–78, 1993.
- [6] P. Domingos. The RISE 2.0 system: A case study in multistrategy learning. Technical Report 95-2, Department of Information and Computer Science, University of California at Irvine, Irvine, CA, 1995.
- [7] P. Domingos. Rule induction and instance-based learning: A unified approach. In *Proceedings of the Fourteenth International Joint Conference on Artificial Intelligence*, Montreal, Canada, 1995. Morgan Kaufmann.
- [8] A. R. Golding and P. S. Rosenbloom. Improving rule-based systems through case-based reasoning. In *Proceedings of the Ninth National Conference on Artificial Intelligence*, pages 22–27, Menlo Park, CA, 1991. American Association for Artificial Intelligence.
- [9] R. Michalski and G. Tecuci, editors. *Machine Learning: A Multistrategy Approach*. Morgan Kaufmann, San Mateo, CA, 1994.
- [10] T. M. Mitchell. Generalization as search. *Artificial Intelligence*, 18:203–226, 1982.
- [11] P. M. Murphy and D. W. Aha. UCI repository of machine learning databases. Machine-readable data repository, Department of Information and Computer Science, University of California at Irvine, Irvine, CA, 1995.
- [12] T. Niblett. Constructing decision trees in noisy domains. In *Proceedings of the Second European Working Session on Learning*, pages 67–78, Bled, Yugoslavia, 1987. Sigma.
- [13] J. R. Quinlan. Simplifying decision trees. *International Journal of Man-Machine Studies*, 27:221–234, 1987.
- [14] J. R. Quinlan. *C4.5: Programs for Machine Learning*. Morgan Kaufmann, San Mateo, CA, 1993.
- [15] J. R. Quinlan. Combining instance-based and model-based learning. In *Proceedings of the Tenth International Conference on Machine Learning*, pages 236–243, Amherst, MA, 1993. Morgan Kaufmann.
- [16] S. Salzberg. A nearest hyperrectangle learning method. *Machine Learning*, 6:251–276, 1991.
- [17] C. Schaffer. Analysis of artificial data sets. In *Proceedings of the Second International Symposium on Artificial Intelligence*, 1989.
- [18] C. Stanfill and D. Waltz. Toward memory-based reasoning. *Communications of the ACM*, 29:1213–1228, 1986.
- [19] D. Wettschereck. A hybrid nearest-neighbor and nearest-hyperrectangle algorithm. In *Proceedings of the Ninth European Conference on Machine Learning*, pages 323–335, Catania, Italy, 1994. Springer-Verlag.
- [20] J. Zhang. A method that combines inductive learning with exemplar-based learning. In *Proceedings of the Second International IEEE Conference on Tools for Artificial Intelligence*, pages 31–37, San Jose, CA, 1990. Institute of Electrical and Electronics Engineers, Computer Society.