

Jmax-pruning: A Facility for the Information Theoretic Pruning of Modular Classification Rules

Frederic Stahl and Max Bramer

*Frederic Stahl, Max Bramer University of Portsmouth, School of Computing,
Buckingham Building, Lion Terrace, PO1 3HE Portsmouth, UK {Frederic.Stahl;
Max.Bramer}@port.ac.uk*

Abstract

The Prism family of algorithms induces modular classification rules in contrast to the Top Down Induction of Decision Trees (TDIDT) approach which induces classification rules in the intermediate form of a tree structure. Both approaches achieve a comparable classification accuracy. However in some cases Prism outperforms TDIDT. For both approaches pre-pruning facilities have been developed in order to prevent the induced classifiers from overfitting on noisy datasets, by cutting rule terms or whole rules or by truncating decision trees according to certain metrics. There have been many pre-pruning mechanisms developed for the TDIDT approach, but for the Prism family the only existing pre-pruning facility is *J-pruning*. *J-pruning* not only works on Prism algorithms but also on TDIDT. Although it has been shown that *J-pruning* produces good results, this work points out that *J-pruning* does not use its full potential. The original *J-pruning* facility is examined and the use of a new pre-pruning facility, called *Jmax-pruning*, is proposed and evaluated empirically. A possible pre-pruning facility for TDIDT based on *Jmax-pruning* is also discussed.

Keywords:

J-pruning, Jmax-pruning, Modular Classification Rule Induction, pre-pruning

1. Introduction

The growing interest in the area of data mining has led to various developments for the induction of classification rules from large data samples

in order to classify previously unseen data instances. Classification rule induction algorithms can be categorised in two different approaches the Top Down Induction of Decision Trees (TDIDT) [1] also known as the ‘divide and conquer’ approach and the ‘separate and conquer’ approach. The ‘divide and conquer’ approach induces classification rules in the intermediate form of a decision tree, whereas the ‘separate and conquer’ approach directly induces ‘*IF... THEN...*’ rules. The ‘divide and conquer’ approach can be traced back to the 1960s [2] and resulted in wide selection of classification systems such as C4.5 and C5.0. The ‘separate and conquer’ approach [3] can also be traced back to the 1960s. Its most notable member is the Prism family of algorithms as a direct competitor to the induction of decision trees.

The original Prism algorithm described in [4] identified the tree structure induced by ‘separate and conquer’ as the major handicap of decision trees which makes them vulnerable to overfitting, especially on noisy datasets. Prism has been shown to produce a similar classification accuracy compared with decision trees and in some cases, even outperforms decision trees. This is particularly the case, if the training data is noisy. There is also some recent interest in parallel versions of Prism algorithms, in order to make them scale better on large datasets. The framework described in [5], the Parallel Modular Classification Rule Inducer (PMCRI) allows to parallelise any member of the Prism family of algorithms.

Nevertheless also Prism, like any classification rule induction algorithm, is prone to overfitting on the training data. For decision trees there is a large variety of pruning algorithms that modify the classifier during or after the induction of the tree in order to make the tree more general and reduce unwanted overfitting [6]. For the Prism family of algorithms there is only one method described in the literature, and that is *J-pruning* [7, 17]. *J-pruning* is based on the *J-measure* [8], a information theoretic measure for quantifying the information content of a rule. Also there is an extension of the PMCRI framework mentioned above, the J-PMCRI framework [9] that incorporates *J-pruning* into PMCRI, not only for quality reasons but also because *J-pruning* reduces the number of rules and rule terms induced and thus the runtime of Prism and PMCRI.

This paper examines the *J-pruning* method described in [7, 17] and proposes a new pruning method *Jmax-pruning*, that aims to exploit the *J-measure* further and thus further to improve Prism’s classification accuracy. The basic Prism approach is described in Section 2 and compared with decision trees, also some particular members of the Prism family are introduced

briefly. *J-pruning* is described and discussed in Section 2.3 and *Jmax-pruning* as a new variation of *J-pruning* is introduced and discussed in Section 3.2 and evaluated in Section 4. The ongoing work is described in Section 5 and discusses J-PrismTCS a version of Prism solely based on the *J-measure* and more importantly proposes the development of a version of *Jmax-pruning* for Decision Tree induction algorithms. Section 6 concludes the paper with a brief summary and discussion of the findings presented.

2. The Prism Family of Algorithms

As mentioned in Section 1, the representation of classification rules differs between the ‘divide and conquer’ and ‘separate and conquer’ approaches. The rule sets generated by the ‘divide and conquer’ approach are in the form of decision trees whereas rules generated by the ‘separate and conquer’ approach are modular. Modular rules do not necessarily fit into a decision tree and normally do not. The rule representation of decision trees is the main drawback of the ‘divide and conquer’ approach, for example rules such as:

$$IF A = 1 AND B = 1 THEN class = x$$

$$IF C = 1 AND D = 1 THEN class = x$$

cannot be represented in a tree structure as they have no attribute in common. Forcing these rules into a tree will require the introduction of additional rule terms that are logically redundant, and thus result in unnecessarily large and confusing trees [4]. This is also known as the replicated subtree problem [10]. Cendrowska illustrates the replicated subtree using the two example rules above in [4]. Cendrowska assumes that the attributes in the two rules above comprise three possible values and both rules predict class x , all remaining classes are labelled y . The simplest tree that can express the two rules is shown in Figure 1 and the total rule set encoded in the tree is:

$$IF A = 1 AND B = 1 THEN Class = x$$

$$IF A = 1 AND B = 2 AND C = 1 AND D = 1 THEN Class = x$$

$$IF A = 1 AND B = 3 AND C = 1 AND D = 1 THEN Class = x$$

$$IF A = 2 AND C = 1 AND D = 1 THEN Class = x$$

$$IF A = 3 AND C = 1 AND D = 1 THEN Class = x$$

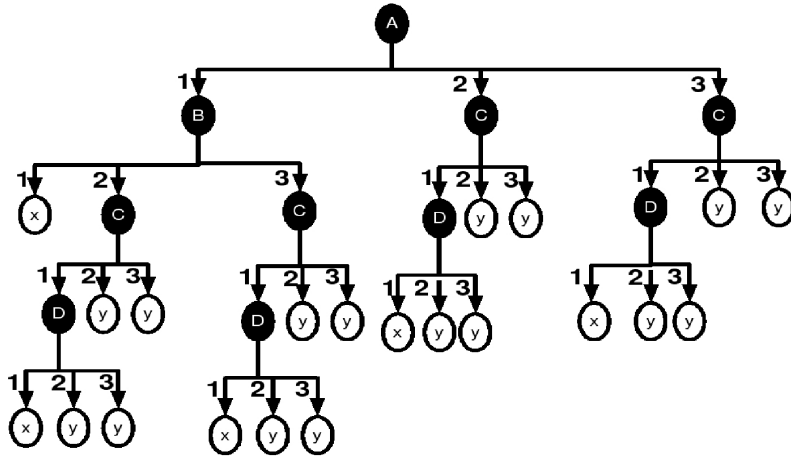


Figure 1: Cendrowska's replicated subtree example.

The fact that modular rules cause trees to grow needlessly complex makes them unsuitable for expert systems as they may require unnecessary expansive tests by the user [4].

'Separate and conquer' algorithms aim to avoid the replicated subtree problem by inducing directly sets of 'modular' rules, avoiding unnecessarily redundant rule terms that are induced just for the representation in a tree structure. The basic 'separate and conquer' approach can be described as follows, where the statement

```
Rule_set rules = new Rule_set();
```

creates a new rule set:

```
Rule_Set rules = new Rule_set();
While Stopping Criterion not satisfied{
    Rule = Learn_Rule;
    Remove all data instances covered from Rule;
    rules.add(rule);
}
```

The *Learn_Rule* procedure generates the best rule for the current subset of the training data where best is defined by a particular heuristic that may vary from algorithm to algorithm. The stopping criterion is also dependent on the algorithm used. After inducing a rule, the rule is added to the rule

set and all instances that are covered by the rule are deleted and a new rule is induced on the remaining training instances.

In Prism each rule is generated for a particular Target Class (TC). The heuristic Prism uses in order to specialise a rule is the probability with which the rule covers the TC in the current subset of the training data. The stopping criterion is fulfilled as soon as there are no training instances left that are associated with the TC.

Cendrowska's original Prism algorithm selects one class as the TC at the beginning and induces all rules for that class. It then selects the next class as TC and resets the whole training data to its original size and induces all rules for the next TC. This is repeated until all classes have been selected as TC. Variations exist such as PrismTC [11] and PrismTCS (Target Class Smallest first) [7]. Both select the TC anew after each rule induced. PrismTC always uses the majority class and PrismTCS uses the minority class. Both variations introduce an order in which the rules are induced, where there is none in the basic Prism approach. However the predictive accuracy of PrismTC cannot compete with that of Prism and PrismTCS (personal communication). PrismTCS does not reset the dataset to its original size and thus is faster than Prism, which produces a high classification accuracy and also sets an order in which the rules should be applied to the test set.

The basic PrismTCS algorithm is outlined below where A_x is a possible attribute value pair and D is the training dataset. The statement

```
Rule_set rules = new Rule_set();
```

creates a new rule set which is a list of rules and the line

```
Rule rule = new Rule(i);
```

creates a new rule with class i as classification. The statement

```
rule.addTerm(Ax);
```

will append a rule term to the rule and

```
rules.add(rule);
```

adds the finished rule to the rule set.

The statements outlined in this paragraph are also used in the pseudo code in Section 2.3 and Section 3.2.

```

    D' = D;
    Rule_set rules = new Rule_set();
Step 1: Find class i that has the fewest instances in
        the training set;
        Rule rule = new Rule(i);
Step 2: Calculate for each Ax p(class = i | Ax);
Step 3: Select the Ax with the maximum p(class = i | Ax);
        rule.addTerm(Ax);
        Delete all instances in D' that do not cover rule;
Step 4: Repeat 2 to 3 for D' until D' only contains
        instances of classification i.
Step 5: rules.add(rule);
        Create a new D' that comprises all instances of
        D except those that are covered by all rules
        induced so far;
Step 6: IF (D' is not empty){
        repeat steps 1 to 6;
    }

```

We will concentrate here on the more popular PrismTCS approach but all techniques and methods outlined here can be applied to any member of the Prism family.

2.1. Dealing with Clashes

A *clash set* is a set of instances in a subset of the training set that are assigned to different classes but cannot be separated further. For example this is inevitable if two or more instances are identical except for their classification. Cendrowska's original Prism algorithm does not take into account that there may be clashes in the training data. However the Inducer software implementations of the Prism algorithms do take clashes into account [12, 11]. What happens in the case of a clash in Inducer is that all instances are treated as if they belong to the TC. [12] mentions that the best approach is to check if the TC is also the majority class. If it is then the rule currently being induced is taken otherwise the rule is discarded. If a clash is encountered and the majority class is not the TC, then the rule is discarded and all instances in the clash set that match the TC are deleted. The reason for manipulating the clash set this way is that if the rule were discarded and the clash set kept then the same rule would be induced all over again and the same clash set would be encountered again.

2.2. Dealing with Continuous Attributes

Continuous attributes are not handled by Cendrowska’s original Prism algorithm. One way to deal with continuous attributes is discretisation of the attribute values prior to the algorithm’s application, for example applying ChiMerge [13] before the application of a Prism algorithm. Bramer’s Inducer software [11] provides implementations of Prism algorithms that deal with continuous attributes. All Prism algorithms implemented in this work deal with continuous attributes the same way as Inducer’s Prism algorithms do. Dealing with continuous attributes can be integrated in step two in the pseudo code above before the calculation of $p(class = i | A_x)$. If A_x is continuous then the training data is sorted for A_x . For example let A_x comprise the following values after sorting, -3, -4.2, 3.5, 5.5 and 10, then the data is scanned for these attribute values in either ascending or descending order. For each attribute value, for example 5.5, two tests $p(class = i | A_x < 5.5)$ and $p(class = i | A_x \geq 5.5)$ are conducted. The one with the largest conditional probability for all the values of the attribute is kept and compared with those conditional probabilities from the remaining attributes.

2.3. J-pruning

As mentioned in the introduction, classifiers are generally pruned to prevent them from overfitting. Pruning methods can be divided into two categories, *pre-pruning* and *post-pruning*. Post-pruning is applied to the classifier after it has been induced whereas pre-pruning is applied during the rule induction process. For Prism algorithms only one pruning method has been developed so far, *J-pruning* [7], a pre-pruning method based on the J-measure [8], a measure for the information content of a rule. *J-pruning* can also be applied to decision tree induction algorithms and has shown good results on both kinds of algorithms [7]. As also mentioned in the introduction, *J-pruning* has found recent popularity in parallel versions of Prism [14, 9, 15], as it reduces the number of rules and rule terms induced considerably and thus increases the computational efficiency.

According to [8] the theoretical average information content of a rule of the form *IF Y = y THEN X = x* can be measured in bits and is denoted by $J(X, Y=y)$.

$$J(X; Y = y) = p(y) \cdot j(X; Y = y) \quad (1)$$

As shown in equation (1) $J(X; Y = y)$ is essentially a product of $p(y)$, the probability with which the left hand side of the rule will occur, and $j(X; Y = y)$ which is called the j-measure (with a lower case j) and measures the goodness-of-fit of a rule. The j-measure, also called the *cross-entropy*, defined in equation (2):

$$j(X; Y = y) = p(x | y) \cdot \log_2\left(\frac{p(x|y)}{p(x)}\right) + (1 - p(x | y)) \cdot \log_2\left(\frac{(1-p(x|y))}{(1-p(x))}\right) \quad (2)$$

For a more detailed description of the J-measure Smyth’s paper [8] is recommended. Bramer’s essential interpretation of the J-measure is that if a rule has a high J-value then it also is likely to have a high predictive accuracy [7]. Hence the J-value is used as an indicator of whether appending further rule terms is likely to improve a rule’s predictive accuracy or lower it due to overfitting. The J-value of the rule may go up or down when appending rule terms, also it may go down and up again. However it is possible to calculate the maximum J-value that the rule with its current terms might maximally achieve if additional terms were added. This upper bound cannot of course be exceeded but its value is not necessarily achievable.

Bramer’s basic *J-pruning* is applied to Prism by calculating the J-value of the rule before the induction of a new rule term and the J-value that the rule would have after a newly induced rule term is appended. If the J-value goes up then the rule term is appended. In the case where the J-value goes down, the rule term is not appended and a test is applied to determine whether the majority class of the instances that are covered by the rule is also the TC. If the majority class is the TC then the rule is truncated and kept and all instances in the current subset of the training set are treated as if they belong to the TC. If the majority class is not the TC, then the rule is discarded and all instances in the clash set that match the TC are deleted, as described in Section 2.1.

A concrete example is used to show how the J-values is used to prune a rule. The example used a dataset extracted from the UCI repository, the soybean dataset [16]. Here we induce rules using our own implementation of PrismTCS with *J-pruning*. The original dataset has been converted to a training set and a test set where the training set comprises 80% of the data instances. The 32nd rule is induced for classification *powdery-mildew*.

First Term

The first rule term induced is (*date = july*), adding up to the intermediate rule:

$$\begin{aligned} & \textit{IF (date = july) THEN CLASS = powdery-mildew} \\ & \text{(J-value = 0.01370)} \end{aligned}$$

Second Term

Now the second rule term is induced, which is in this case (*crop-hist = same-lst-yr*). Now the J-value of the rule is calculated if the second rule terms (indicated in bold letters) were appended.

$$\begin{aligned} & \textit{IF (date = july) \textbf{AND (crop-hist = same-lst-yr)} THEN CLASS =} \\ & \textit{powdery-mildew} \\ & \text{(J-value = 0.01157)} \end{aligned}$$

Now it can be seen that appending the second rule term would lower the J-value of the rule by 0.00213 (0.01370-0.01157), hence the second rule term is not appended leaving rule:

$$\textit{IF (date = july) THEN CLASS = powdery-mildew}$$

This is checked to ensure that the target class *powdery-mildew* is the majority class of the instances under consideration. This is true in this case, hence the rule is taken into the rule set. If the majority class were not *powdery-mildew* the rule would be discarded and clash handling as described in Section 2.1 would be applied.

The basic PrismTCS algorithm incorporating J-pruning is outlined below where A_x is a possible attribute value pair, D is the training dataset and C the clash set which is empty at the beginning of the algorithm's execution. Some of the statements used in this pseudo code are outlined for the pseudocode of PrismTCS above. The Basic J-pruning technique can be found in Step 4:

```

initialise:
    Rule_set rules = new Rule_set;
    D' = D;
Step 1: Find class i that has the fewest instances in the
        training set;
        double bestJValue = 0;
        Rule rule = new Rule(i);
Step 2: Calculate for each Ax p(class = i | Ax);
Step 3: Select the Ax with the maximum p(class = i | Ax)
        and calculate the rule's J-value (tempJValue) if
        Ax would be appended to the rule;
Step 4: IF(tempJValue > bestJValue){
        rule.addTerm(Ax);
        bestJValue = tempJValue;
        Delete all instances in D' that do not cover rule;
    } ELSE{
        IF((D' has i as majority class) AND
            (bestJValue > 0)){
            rules.add(rule);
            GO TO Step 6;
        }ELSE{
            C = C + (instances in D' covered by rule
                and Ax and match class i);
            D = D - C;
            discard the current rule;
            GO TO Step 6;
        }
    }
}
Step 5: Repeat 2 to 4 for D' until D' only contains instances
        of classification i. The induced rule is then a
        conjunction of all selected Ax and i;
Step 6: Create a new D' that comprises all instances of
        D except those that are covered by all rules
        induced so far;
Step 7: IF (D' is not empty){
        repeat steps 1 to 6;
    }
}

```

3. Variation of J-pruning

In general there is very little work on pruning methods for the Prism family of algorithms. Bramer’s *J-pruning* in the Inducer software seems to be the only pruning facility developed for Prism algorithms. This Section critiques the initial *J-pruning* facility and outlines *Jmax-pruning*, a variation that makes further use of the J-measure.

3.1. Critique of J-pruning

Even though *J-pruning* described in Section 2.3 achieves good results regarding the overfitting of Prism, it does not exploit the J-measure to its full potential. The reason is that even if the new rule term decreases the J-value, it is possible that the J-value increases again when adding further rule terms [8]. If the rule is truncated as soon as the J-value is decreased it may result in the opposite of overfitting, an over generalised rule with a lower predictive accuracy. The relatively good results for *J-pruning* achieved in [7] could be explained by the assumption that it does not happen very often that the J-value decreases and then increases again. However, how often this happens will be examined empirically in Section 4.

3.2. Jmax-pruning

According to [8], an upper bound for the J-measure for a rule can be calculated using equation (3):

$$J_{max} = p(y) \cdot \max\left\{ p(x | y) \cdot \log_2\left(\frac{1}{p(x)}\right), (1 - p(x | y)) \cdot \log_2\left(\frac{1}{1-p(x)}\right) \right\} \quad (3)$$

If the actual J-value of the rule currently being generated term by term matches the maximum possible J-value (J_{max}) it is an absolute signal to stop the induction of further rule terms.

A concrete example is used to show how the J-values of a rule can develop. The example used a dataset extracted from the UCI repository, the soybean dataset [16]. Here we induce rules using our own implementation of PrismTCS without any *J-pruning*. The original dataset has been converted to a training set and a test set where the training set comprises 80% of the data instances. The 39th rule induced is:

IF (temp = norm) AND (same-lst-sev-yrs = whole-field) AND (crop-hist = same-lst-two-yrs) THEN CLASS = frog-eye-leaf-spot

This is a perfectly reasonable rule with a J-value of 0.00578. However looking at the development of the J-values after each rule term appended draws a different picture:

First Term

IF (temp = norm) THEN CLASS = frog-eye-leaf-spot

(J-value = 0.00113, J_{max} = 0.02315)

Here the rule has J-value of 0.00113 after the first rule term has been appended. The J-value for the complete rule (0.00578) is larger than the current J-value, which is to be expected as the rule is not fully specialised yet on the TC.

Second Term

IF (temp = norm) AND (same-lst-sev-yrs = whole-field) THEN CLASS = frog-eye-leaf-spot

(J-value = 0.00032, J_{max} = 0.01157)

Now the J-value is decreased to 0.00032 and J_{max} to 0.01157. Here *J-pruning* as described in Section 2.3 would stop inducing further rule terms, the finished rule would be:

IF (temp = norm) THEN CLASS = frog-eye-leaf-spot

with a J-value of 0.00113. However looking at the value of J_{max} , after the second rule term has been appended, it can be seen that it is still higher than the previous J-value for appending the first rule term. Thus it is still possible that the J-value may increase again above the so far highest J-value of 0.00113. Inducing the next rule term leads to:

Third Term

IF (temp = norm) AND (same-lst-sev-yrs = whole-field) AND (crop-hist = same-lst-two-yrs) THEN CLASS = frog-eye-leaf-spot

$$(J\text{-value} = 0.00578, J_{max} = 0.00578)$$

In this case the rule was finished after the appending of the third rule term as it only covered examples of the TC. However, the interesting part is that the J-value increased again by appending the third rule term and is in fact the highest J-value obtained. Bramer’s original method would have truncated the rule too early leading to an overall average information content of 0.00113 instead of 0.00578. The J-value and the J_{max} value are rounded to five digits after the decimal point and appear identical but are actually slightly different. Looking at more digits the values are in fact for the J-value 0.005787394940853119 and for the J_{max} 0.005787395266794598. In this case no further rule terms can be added to the left-hand side of the rule as the current subset of the training set only contains instances of the TC, but if this were not the case it would still not be worthwhile to add additional terms as the J-value is close to J_{max} .

Overall this observation strongly suggests that pruning the rule as soon as the J-value decreases does not fully exploit the J-measure’s potential. This work suggests that *J-pruning* could be improved by inducing the maximum number of possible rule terms until the current subset cannot be broken down further or the actual J-value is equal to or within a few percent of J_{max} . As a rule is generated all the terms are labelled with the actual J-value of the rule after appending this particular rule term. The partial rule for which the largest rule’s J-value was calculated would then be identified and all rule terms appended afterwards truncated, with clash handling as described in Section 2.1 invoked for the truncated rule. We call this new pre-pruning method *Jmax-pruning*.

The basic PrismTCS algorithm incorporating Jmax-pruning is outlined below where A_x is a possible attribute value pair, D is the training dataset and C the clash set. Some of the statements used in this pseudo code are outlined in the PrismTCS pseudo code in Section 2. However the statement

```
rule.removeTerm(termIndex);
```

is new, it removes a term at the specified index. The *Jmax-pruning* pseudo code can be found in Step 4:

```
initialise:
    D' = D;
```

```

Rule_Set rules;
Step 1: D'' = D';
Find class i that has the fewest instances in D'';
Rule rule = new Rule(i);
int termCounter = 0;
double bestJValue = 0;
int bestTermIndex = 0;
Step 2: Calculate for each Ax p(class = i | Ax);
Step 3: Select the Ax with the maximum p(class = i | Ax);
rule.addTerm(termCounter, Ax);
Delete all instances in D'' that do not cover rule;
double tempJValue = new J-value for rule;
IF(tempJValue > bestJValue){
    bestJValue = tempJValue;
    bestTermIndex = termCounter;
}
termCounter++;
delete all instances in D'' that are not covered by
the selected Ax;
Step 4: IF(D'' only contains instances of classification i){
    IF(termCounter == bestTermIndex){
        rules.add(rule);
    }ELSE{
        FOR(termCounter; termCounter > bestTermIndex;
            termCounter--){
            rule.removeTerm(termCounter);
        }
        IF(instances in D' covered by rule having
            i as majority class){
            rules.add(rule);
        }ELSE{
            C = C + (instances in D' covered by rule not
                having i as majority class);
            delete rule;
        }
    }
}ELSE{
    Repeat Steps 2 to 4;
}
Step 6: D' = D - C;

```

```

Delete all instances in D' that cover rules;
Step 7: IF (D' is not empty){
    repeat steps 1 to 6;
}

```

An interesting fact to note is that *Jmax-pruning* can be seen as a hybrid between pre- and post-pruning. That is because rules are firstly fully induced and then pruned which applies to post pruning, however the fact that each induced rule is ‘post-pruned’ before the next rule is induced can be pre-pruning. Contrary to *Jmax-pruning*, *J-pruning* is completely based on pre-pruning as rules are not fully induced before they are pruned.

4. Evaluation of Jmax-pruning

The datasets used have been retrieved from the UCI repository [16]. Each dataset is divided into a test set holding 20% of the instances and a training set holding the remaining 80% of the instances.

Table 1 shows the number of rules induced per training set and the achieved accuracy on the test set using PrismTCS with *J-pruning* as described in Section 2.3 and *Jmax-pruning* as proposed in Section 3.2.

What is also listed in Table 1 as ‘J-value recovers’, is the number of times the J-value decreased and eventually increased again when first fully expanding the rule and then pruning it using *Jmax-pruning*. Using the original *J-pruning* as described in Section 2.3 would not detect these J-value recoveries and lead to a rule with a lower J-value and thus lower information content than it could possibly achieve.

What can be seen is that in all cases *Jmax-pruning* performs either better than or produces the same accuracy as *J-pruning*. In fact seven times *Jmax-pruning* produced a better result than *J-pruning* and nine times it produced the same accuracy as *J-pruning*. Taking a closer look in the rule sets that have been produced in the nine cases for which the accuracies for both pruning methods are the same revealed that identical rule sets were produced in seven out of these nine cases. The two exceptions are the ‘Car Evaluation’ and ‘ecoli’ datasets, however in these two exceptions the classification accuracy was the same using *Jmax-pruning* or *J-pruning*. In the cases where there are identical classifiers there were no J-value recoveries present. In Section 3.1 we stated that the good performance of *J-pruning* [7], despite its tendency to generalisation, can be explained by the fact that there are not many J-value recoveries in the datasets and thus the tendency to over generalisation

Table 1: Comparison of *J-pruning* and *Jmax-pruning* on PrismTCS.

Dataset	Number of Rules	Accuracy (%)	Number of Rules	Accuracy (%)	J-value recovers
	J-Pruning		J-max Pruning		
monk1	4	79	12	86	4
monk3	3	98	3	98	0
vote	3	94	3	94	0
genetics	8	70	8	70	0
22 contact lenses	4	95	4	95	0
breast cancer	24	96	24	96	0
soybean	39	88	43	89	4
australian credit	20	89	20	89	0
diabetes	29	75	31	76	1
crx	18	83	18	83	0
segmentation	83	79	86	82	2
ecoli	23	78	26	78	3
Balance Scale	10	72	36	74	21
Car Evaluation	4	76	4	76	1
Contraceptive Method Choice	19	44	28	45	8
Optical Recognition of handwritten Digits	456	57	467	58	6

is low. Looking into the last column of table 1 we can see the number of J-value recoveries. In seven cases there are none, thus there is no potential for over generalisation by using *J-pruning* and for the remaining datasets there is only a very small number of J-value recoveries with the exception of the ‘Balanced Scale’ dataset for which a 2% higher accuracy has been retrieved by using *Jmax-pruning* compared with *J-pruning*.

Loosely speaking, if there are no J-value recoveries present, then Prism algorithms with *Jmax-pruning* will produce identical classifiers to Prism algorithms with *J-pruning*. However, if there are J-value recoveries, it is likely that Prism algorithms with *Jmax-pruning* will produce classifiers that achieve a better accuracy than Prism algorithms with *J-pruning*.

What can also be read from Table 1 is the number of rules induced. In all cases in which both pruning methods produced the same accuracy, the classifiers and thus the number of rules were identical. However in the

cases where the J-value recovered, then the number of rules induced with *Jmax-pruning* was the same or larger than the number of rules induced with *J-pruning*. This can be explained by the fact that in the case of a J-value recovery the rule gets specialised further than with normal *J-pruning* by adding more rule terms while still avoiding overfitting. Adding more rule terms results in the rule covering fewer training instances from the current subset. This in turn results in that before the next iteration for the next rule fewer instances are deleted from the training set, which potentially generates more rules, assuming that the larger the number of training instances the more rules are generated.

5. Ongoing Work

5.1. *J-PrismTCS*

Another possible variation of PrismTCS that is currently being implemented is a version that is solely based on the J-measure. Rule terms would be induced by generating all possible categorical and continuous rule terms and selecting the one that results in the highest J-value for the current rule instead of selecting the one with the largest conditional probability. Again the same stopping criterion as for standard PrismTCS could be used, which is that all instances of the current subset of the training set belong to the same class. We call this variation of PrismTCS, *J-PrismTCS*.

5.2. *Jmax-pruning for TDIDT*

J-pruning has also been integrated into the TDIDT approach as a pre-pruning facility and achieved a higher classification accuracy than TDIDT without *J-pruning* [7]. Encouraged by the good results outlined in Section 4, which were achieved with *Jmax-pruning* in PrismTCS, we are currently developing a version of *Jmax-pruning* for TDIDT algorithms. The following pseudo code describes the basic TDIDT algorithm.

```

IF    All instances in the training set belong to the
      same class
THEN  return value of this class
ELSE  (a) Select attribute A to split on
      (b) Divide instances in the training set
           into subsets, one for each value of A.
      (c) Return a tree with a branch for each non

```

empty subset, each branch having a decendent subtree or a class value produced by applying the algorithm recursively

The basic approach of *J-pruning* in TDIDT is to prune a branch in the tree as soon as a node is generated at which the J-value is less than at its parent node [7]. However performing the *J-pruning* is more complicated than for Prism algorithms as illustrated in the example below.

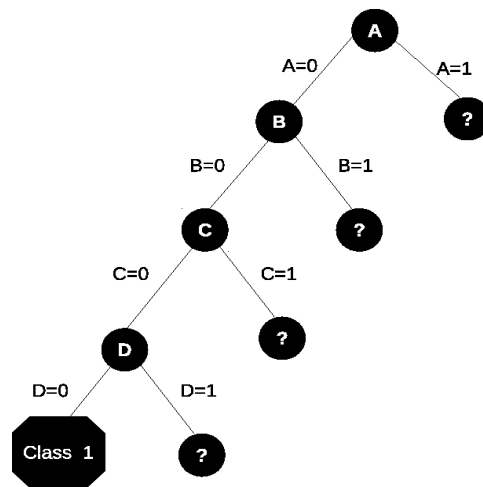


Figure 2: Example of a decision tree for *J-pruning*.

Figure 2 illustrates a possible tree which is used to explain *J-pruning* and *Jmax-pruning* for TDIDT. The nodes labelled with ‘?’ are placeholders for possible subtrees. Now assuming that a depth first approach is used and the current node being expanded is node ‘D’. In this case *J-pruning* would take the incomplete rule

(1) *IF (A=0) AND (B=0) AND (C=0)...*

the complete rule

(2) *IF (A=0) AND (B=0) AND (C=0) AND (D=0) THEN class = 1*

and the possible incomplete rule

(3) IF (A=0) AND (B=0) AND (C=0) AND (D=1)...

into account. Rule (2) is completed as all instances correspond to the same classification which is (class = 1). However instances covered by incomplete rule (3) correspond in this case to more than one classification.

J-pruning now compares the J-values of the incomplete rules (1) and (3). If the J-value of rule (3) is less than the J-value of rule (1) then rule (3) is completed by assigning it to the majority class of the corresponding instances. The complication is that the calculation of the J-value of a rule requires us to know its right-hand side. In the case of a complete (non-truncated) rule, such as rule (2), this is straightforward, but how can the J-value be calculated for an incomplete rule?

The method described by Bramer [7] is to imagine all possible alternative ways of completing the incomplete rule with right-hand sides class=1, class=2 etc., calculate the J-value of each such (completed) rule and take the largest of the values as the estimate of the J-value of the incomplete rule.

In a similar way to J-pruning for Prism algorithms, J-pruning for TDIDT in its current form does not necessarily exploit the full potential of the J-measure as again it is possible that if rule (3) were not truncated at the node labelled ‘?’ but expanded to complete the decision tree in the usual way the J-value for some or possibly all of the resulting complete branches might be at least as high as the J-value at node D.

Applying the idea of Jmax-pruning rather than J-pruning to TDIDT may increase the classification accuracy. This could be done by developing the complete decision tree and labelling each internal node with a J-value estimated as described above. Each branch (corresponding to a completed rule) can then be truncated at the node that gives the highest of the estimated J-values, in a similar way to the method described in Section 3.2, with each truncated rule assigned to the majority class for the corresponding set of instances.

This method appears attractive but there is a possible problem. Using the example from Figure 2 and assuming that the estimated J-value of rule (1) is greater than the estimated J-value of rule (3) and that the majority class of the instances at node D is ‘1’, then rule (1) would be truncated at node D and rule (3) would cease to exist, giving two completed rules in this subtree:

(1) IF (A=0) AND (B=0) AND (C=0) THEN class = 1

and

(2) IF (A=0) AND (B=0) AND (C=0) AND (D=0) THEN class = 1

Both rules are illustrated in Figure 3. Rule (2) is now redundant. It is just a special case of rule (1), with the same classification, and can be discarded.

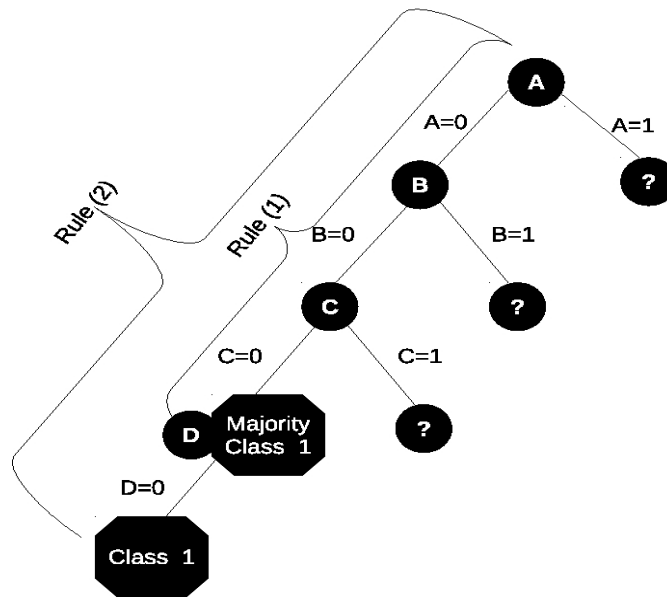


Figure 3: Example of a decision tree with a redundant rule

Now suppose instead that in the above the majority class of the instances at node (1) were '2' (rather than '1'). In this case a different picture would emerge, with rules (1) and (2) having different classifications. How likely this situation is to occur in practice and how it should best be handled if it does both remain to be determined.

6. Conclusions

Section 2 discussed the replicated subtree problem introduced due to the representation of rules in the form of decision trees. The Prism family of algorithms has been introduced as an alternative approach to TDIDT

that does induce modular classification rules that do not necessarily fit into a tree structure. The Prism family of algorithms was highlighted and *J-pruning*, a pre-pruning facility for Prism algorithms based on the J-measure, which describes the information content of a rule, was introduced. Section 3 criticised *J-pruning* as it does not fully exploit the potential of the J-measure. The J-value of a rule may go up or down when rule terms are appended to the rule. *J-pruning* truncates a rule as soon as the J-value decreases even if it may recover (increase again). The proposed *Jmax-pruning* exploits the possibility of a J-value recovery and achieves in some cases, examined in Section 4, better results compared with *J-pruning*, but in every case examined *Jmax-pruning* achieved at least the same or a higher classification accuracy compared with *J-pruning*.

The ongoing work comprises the development of J-PrismTCS, a version of PrismTCS that is solely based on the J-measure, by using it also as a rule term selection metric as discussed in Section 5.1. Furthermore the ongoing work comprises the development of a TDIDT algorithm that incorporates *Jmax-pruning* as discussed in Section 5.2.

References

- [1] J. R. Quinlan, C4.5: programs for machine learning, Morgan Kaufmann, 1993.
- [2] E. B. Hunt, P. J. Stone, J. Marin, Experiments in induction, Academic Press, New York, 1966.
- [3] R. S. Michalski, On the Quasi-Minimal solution of the general covering problem, in: Proceedings of the Fifth International Symposium on Information Processing, Bled, Yugoslavia, pp. 125–128.
- [4] J. Cendrowska, PRISM: an algorithm for inducing modular rules, International Journal of Man-Machine Studies 27 (1987) 349–370.
- [5] F. T. Stahl, M. A. Bramer, M. Adda, PMCRI: A parallel modular classification rule induction framework, in: MLDM, Springer, 2009, pp. 148–162.
- [6] F. Esposito, D. Malerba, G. Semeraro, A comparative analysis of methods for pruning decision trees, IEEE Transactions on Pattern Analysis and Machine Intelligence 19 (1997) 476–491.

- [7] M. A. Bramer, An information-theoretic approach to the pre-pruning of classification rules, in: B. N. M Musen, R. Studer (Eds.), *Intelligent Information Processing*, Kluwer, 2002, pp. 201–212.
- [8] P. Smyth, R. M. Goodman, An information theoretic approach to rule induction from databases, *Transactions on Knowledge and Data Engineering* 4 (1992) 301–316.
- [9] F. Stahl, M. Bramer, M. Adda, J-PMCRI: A methodology for inducing pre-pruned modular classification rules, in: *Artificial Intelligence in Theory and Practice III*, Springer, Brisbane, 2010, pp. 47–56.
- [10] I. H. Witten, F. Eibe, *Data Mining: Practical Machine Learning Tools and Techniques with Java Implementations*, Morgan Kaufmann, 1999.
- [11] M. A. Bramer, Inducer: a public domain workbench for data mining, *International Journal of Systems Science* 36 (2005) 909–919.
- [12] M. A. Bramer, Automatic induction of classification rules from examples using N-Prism, in: *Research and Development in Intelligent Systems XVI*, Springer-Verlag, Cambridge, 2000, pp. 99–121.
- [13] R. Kerber, Chimerge: Discretization of numeric attributes, in: *AAAI*, pp. 123–128.
- [14] F. T. Stahl, *Parallel Rule Induction*, Ph.D. thesis, University of Portsmouth, 2009.
- [15] F. T. Stahl, M. Bramer, M. Adda, Parallel rule induction with information theoretic pre-pruning, in: *SGAI Conf.*, pp. 151–164.
- [16] C. L. Blake, C. J. Merz, *UCI repository of machine learning databases*, Technical Report, University of California, Irvine, Department of Information and Computer Sciences, 1998.
- [17] M. A. Bramer, Using J-Pruning to Reduce Overfitting in Classification Trees, *Knowledge-Based Systems* 15 (2002) 301–308.