

# Attribute Generation Based on Association Rules

Masahiro Terabe<sup>1</sup>, Takashi Washio<sup>2</sup>, Hiroshi Motoda<sup>2</sup>,  
Osamu Katai<sup>3</sup> and Tetsuo Sawaragi<sup>4</sup>

<sup>1</sup>Mitsubishi Research Institute Inc., Tokyo, Japan

<sup>2</sup>Institute of Scientific and Industrial Research, Osaka University, Osaka, Japan

<sup>3</sup>Graduate School of Informatics, Kyoto University, Kyoto, Japan

<sup>4</sup>Graduate School of Engineering, Kyoto University, Kyoto, Japan

**Abstract.** A decision tree is considered to be appropriate (1) if the tree can classify the unseen data accurately, and (2) if the size of the tree is small. One of the approaches to induce such a good decision tree is to add new attributes and their values to enhance the expressiveness of the training data at the data pre-processing stage. There are many existing methods for attribute extraction and construction, but constructing new attributes is still an art. These methods are very time consuming, and some of them need a priori knowledge of the data domain. They are not suitable for data mining dealing with large volumes of data. We propose a novel approach that the knowledge on attributes relevant to the class is extracted as association rules from the training data. The new attributes and the values are generated from the association rules among the originally given attributes. We elaborate on the method and investigate its feature. The effectiveness of our approach is demonstrated through some experiments.

**Keywords:** Association rules; Attribute generation; Data mining; Data pre-processing; Decision tree

---

## 1. Introduction

Data mining is becoming a key technique in discovering meaningful patterns and rules from large amounts of data (Berry and Linoff, 1997). It is often used in the fields of business such as marketing and customer support operations. However, one of the major difficulties in the practical approach is scarceness of experts in the application domain and good data-mining tools.

In the data-mining process, a decision tree is often used as knowledge representation of the mined result, because it is easy to understand for human analysts.

---

*Received 6 December 1999*

*Revised 28 October 2000*

*Accepted 9 March 2001*

In general, the appropriateness of the decision tree is evaluated from the next two criteria.

- *Size*: The size of decision tree is evaluated by the number of nodes in the tree. A smaller decision tree is easier to understand. It is also known that a smaller decision tree leads to avoiding overfitting to the data.
- *Prediction accuracy*: The decision tree predicting the correct class of a new instance with higher accuracy is desired.

However, inducing an appropriate decision tree is often difficult. In many cases, the original attributes are not expressive enough. Further, some of them are irrelevant or redundant. Feature selection removes irrelevant and redundant attributes (Liu and Motoda, 1998). Feature extraction and construction create new attributes (Bloedorn and Michalski, 1998; Lavrač et al., 1998) and add them to the original training data. As a result, the description of the original training data is enriched. The decision tree induced from the pre-processed training data can be better than that induced from the original training data. Constructing new attributes without domain knowledge is computationally very expensive. Preparing appropriate construction operators and applying them in the right order is still an art. These processes are very time consuming and some of them need a priori knowledge of the data domain. These features are not suitable for data mining. When the method is applied to data mining, the computational cost of the method should be reasonably small even if the data size is large. Furthermore, domain knowledge is insufficient at the beginning of the data-mining process, because the human analyst may not be an expert in the data domain.

In this paper, we propose a novel method of attributes generation, which has suitable properties as a data-mining method. Our proposed method generates new attributes based on association rules among the original attributes. The rules are extracted from the original training data automatically.

The proposed method has the following properties:

1. It does *not need any a priori knowledge of the attributes and their association*. The knowledge is extracted from the original training data automatically.
2. It adopts the Apriori algorithm to generate *attributes association rules* so that the rules are extracted with *reasonable computational cost even if the data size is large*.

These properties are very advantageous compared to the aforementioned traditional attribute generation methods. The first property makes it easy to use the method from the beginning of the data-mining process when domain knowledge is insufficient. The second enables to apply the method to large scale data.

The paper is organized as follows. In Section 2, we briefly explain association rules and Apriori algorithm. In Section 3, we propose the novel attribute generation method based on the Apriori algorithm. In Section 4, we investigate the performance of our proposed method. In Section 5, we discuss the characteristics of the proposed method based on experimental results.

## 2. Association Rules

The Apriori algorithm (Agrawal and Srikant, 1994) extracts a co-occurrence pattern of items from the instances in the following form of an association rule.

$$R : \mathbf{B} \Rightarrow \mathbf{H}, \quad (1)$$

where

$\mathbf{B}$  : Body which is the condition part of the association rule, and

$\mathbf{H}$  : Head which is the conclusion part of the association rule.

Both 'Body' and 'Head' are a set of items. This association rule means that 'If an instance includes all the items in *Body*, then the instance also includes all the items in *Head* in many cases.' In other words, the association rule indicates a co-occurrence pattern of item sets in the data.

Traditional algorithms need much computation time to extract association rules. The Apriori algorithm proposed by Agrawal succeeded in reducing the search space efficiently so that the computation time is much smaller than traditional methods even if the data size is large.

In the Apriori algorithm, the candidates of the association rule are evaluated by two indices, i.e., 'support value' and 'confidence value'. The support value  $sup(R)$  of the association rule  $R$  is defined as follows:

$$sup(R : \mathbf{B} \Rightarrow \mathbf{H}) = \frac{n(\mathbf{B} \cup \mathbf{H})}{N}, \quad (2)$$

where

$n(\mathbf{B} \cup \mathbf{H})$  : number of instances which include all items in both  $\mathbf{B}$  and  $\mathbf{H}$ , and

$N$  : total number of data.

The support value indicates the ratio of the number of instances including all items appearing in the rule to all instances. Therefore, the rule covers larger portion of the data if the support value of the rule is higher.

The confidence value of the association rule  $conf(R)$  is defined as follows:

$$conf(R : \mathbf{B} \Rightarrow \mathbf{H}) = \frac{n(\mathbf{B} \cup \mathbf{H})}{n(\mathbf{B})}, \quad (3)$$

where

$n(\mathbf{B})$  : number of instances which include all items in  $\mathbf{B}$ .

If the confidence value is higher, the association is more plausible.

In the Apriori algorithm, both the minimum support value and the minimum confidence value are given as threshold parameters. In the association rule extraction process, the rules which do not satisfy these threshold conditions are removed from the set of candidates. The search uses the monotonicity of support values; i.e., if the support of an item set is below the threshold, its super set is pruned from the search. If these minimum thresholds given are decreased, more candidates for association rules are generated and evaluated in the algorithm. As a result, the computation time increases, though a more complete set of good association rules would be extracted. On setting these minimum thresholds, the trade-off between the cost of computation time and the rule covering good association should be taken into account carefully.

### 3. Proposed Attribute Generation Method

First, the proposed method extracts *Attributes Association Rules (AARs)* as the basic knowledge of associations among attributes from the original training data prepared for decision tree induction. Next, the proposed method generates some new attributes based on the extracted AARs, and adds them to the original training data. Therefore, the data description is extended by this data pre-processing. After this data pre-processing to the original training data, a decision tree is induced by the standard decision tree algorithm.

The AARs represent associations among attributes and the class. Our proposed method consists of the following steps:

- Step 1:** Description of training data for a decision tree is transformed to the transaction format.
- Step 2:** Apriori algorithm extracts AARs from the transaction data.
- Step 3:** Some candidates of new attributes are generated based on AARs.
- Step 4:** Degree of contribution of the new attribute candidates to identifying the class is evaluated.
- Step 5:** The new attribute candidates that satisfy a criterion (explained later) are added.

The details of each step are explained in the following subsections.

#### 3.1. Step 1: Data Description Transformation

First, the proposed method transforms the data description of the original training data to the transaction format.

The training data is described as a set of instance data, *train\_data*. An instance in *train\_data*, *datum* is represented as follows:

$$datum = \{v_{i,j} | \forall i \in M, \exists j \in N_i\} \cup \{v_{c,j} | \exists j \in N_c\}, \quad (4)$$

where

$$\begin{aligned} M &= \{1, \dots, m\}, \\ N_i &= \{1, \dots, n_i\}, \\ N_c &= \{1, \dots, n_c\} \end{aligned}$$

where  $v_{i,j}$  is a value of an attribute  $a_i$ ,  $m$  the number of attributes,  $n_i$  the total number of values of  $a_i$ , and  $n_c$  the total number of values of the class  $c$ .

In the proposed method, each pair of attribute ' $a_i$ ' and its value ' $v_{i,j}$ ' is transformed to an item in the form of  $item_i = \langle a_i, v_{i,j} \rangle$  and  $item_{m+1} = \langle c, v_{c,j} \rangle$ . Under this transformation, *datum* in equation (4) becomes a transaction as follows:

$$trans = \{item_1, \dots, item_i, \dots, item_{m+1}\} \quad (5)$$

The transaction data, *trans\_data*, is a set of transactions.

#### 3.2. Step 2: Extraction of Attributes Association Rules

Next, the proposed method extracts AARs that represent the association among pairs of attributes and its attribute value.

The attribute value that has strong association with the class value is useful in predicting the class. Accordingly, the association between the attributes and the class is a measure of the goodness of the attribute. For that reason, we extract AARs, which satisfy the following two conditions:

- The condition part includes only an item set consisting of attributes.
- The conclusion part includes only an item representing the class.

An AAR  $R$  that satisfies these conditions is described as follows:

$$R : \text{if } \bigcup_{\forall i \in M_s} \{ \langle a_i, v_{i,j} \rangle \in ITEM_{a_i} | \exists j \in N_i \} \\ \text{then } \{ \langle c, v_{c,j} \rangle \in ITEM_c | \exists j \in N_c \}, \quad (6)$$

where

$$M_s \subseteq M,$$

$$ITEM_{a_i} = \bigcup_{\forall \langle a_i, v_{i,j} \rangle \in \forall trans \in trans\_data} \{ \langle a_i, v_{i,j} \rangle \},$$

$$ITEM_c = \bigcup_{\forall \langle c, v_{c,j} \rangle \in \forall trans \in trans\_data} \{ \langle c, v_{c,j} \rangle \}$$

This represents a fact that the instance involving these pairs of the attribute and its value in the condition part is concluded with the class  $v_{c,j}$  with high confidence. Consequently, a new composed attribute which is the collection of the attributes and their values in these pairs is expected to be useful in predicting the class of an instance.

### 3.3. Step 3: Generating New Attribute Candidates

Let  $\mathcal{R}$  be a set of all AARs extracted in Step 2. Each AAR  $R \in \mathcal{R}$  is a basic unit of this attribute generation algorithm. Let  $\mathbf{B}$  is the body of  $R$ ,  $\mathbf{H}$  is the head of  $R$ ,  $A(\mathbf{B}) = \{ a_i | \forall \langle a_i, v_{i,j} \rangle \in \mathbf{B} \}$ , and  $C(\mathbf{H}) = \{ v_{c,j} | \langle c, v_{c,j} \rangle \in \mathbf{H} \}$ . Define a partition  $\mathcal{P}$  such that each element of  $\mathcal{P}$  is  $P_q = \{ R | \forall A(\mathbf{B}); A(\mathbf{B})\text{s are mutually identical in } \mathcal{R} \}$ .

A new attribute candidate is characterized by the following quadruple  $AN_q$  for each  $P_q$ :

$$AN_q = \langle \mathcal{A}_q, \mathcal{V}_q, \mathcal{S}_q, \mathcal{C}_q \rangle, \quad (7)$$

where

$$\mathcal{A}_q : A(\mathbf{B}) \text{ where } \mathbf{B} \text{ is the body of } R \text{ in } P_q,$$

$$\mathcal{V}_q = \{ V_{q_0}, V_{q_1}, \dots, V_{q_{|P_q|}} \},$$

$$\mathcal{S}_q = \{ sup(R) | \forall R \in P_q \},$$

$$\mathcal{C}_q = \{ conf(R) | \forall R \in P_q \}$$

where

$$V_{q_0} = \bigwedge_{\forall R \in P_q} \neg true(\mathbf{B})$$

$$V_{q_k} \in \{ true(\mathbf{B}) | \forall R \in P_q \} \quad (k = 1, \dots, |P_q|)$$

The attributes in the Body  $A(\mathbf{B})$  where  $\mathbf{B}$  is the body of AARs  $R$  in  $P_q$  are

**Table 1.** New attribute generation algorithm: *NewAttributeCand*


---

**A set of AARs extracted in Step 2:**  $\mathcal{R}$   
**An AAR:**  $R_p(\in \mathcal{R}) : B_p \Rightarrow H_p$ ,  
 where  
 $B_p = \bigcup_{v_i \in M_{sp}} \{ \langle a_i, v_{i,j} \rangle \in ITEM_{a_i} | \exists j \in N_i \}$ ; /\*condition part\*/  
 $H_p = \{ \langle c, v_{c,j} \rangle \in ITEM_c | \exists j \in N_c \}$ ; /\*conclusion part\*/  
 $sup(R_p)$ ; /\*support\*/  $conf(R_p)$ ; /\*confidence\*/  
 $A(B_p) = \{ a_i | \forall \langle a_i, v_{i,j} \rangle \in B_p \}$ ;  
**A set of new attribute candidate information:**  $\mathcal{AN}$   
**A new attribute candidate information:**  $AN_q = \langle \mathcal{A}_q, \mathcal{V}_q, \mathcal{S}_q, \mathcal{C}_q \rangle \in \mathcal{AN}$ ,  
 where  
 $\mathcal{A}_q$ ; /\* new attribute candidate\*/  
 $\mathcal{V}_q$ ; /\* a set of attribute values\*/  
 $\mathcal{S}_q$ ; /\* a set of support value of each original AAR \*/  
 $\mathcal{C}_q$ ; /\* a set of confidence value of each original AAR \*/

**Algorithm:**  
*NewAttributeCand*( $\mathcal{R}, \mathcal{AN}$ ) {  
 /\* Generate a partition  $\mathcal{P}$  such that each element of \*/  
 /\*  $\mathcal{P}$  is a partition such that the element \*/  
 /\*  $P_q = \{ R_p | \forall A(B_p), A(B_p) \text{ s are mutually identical in } \mathcal{R} \}$  \*/  
 $q = 1; P_1 = \{ R_1 \}; \mathcal{P} \leftarrow \{ P_1 \}$ ;  
**for**( $p = 2; p \leq |\mathcal{R}|; p++$ ) {  
 /\*  $A(B_p)$  is identical with a  $A(B)$  where  $B$  is the body of  $R \in P_q \in \mathcal{P}$  \*/  
**if**( $\exists P_q \in \mathcal{P} | A(B_p) == A(B); B$  is the body of  $R \in P_q$ ) {  
 /\* Add  $R_p$  to  $P_q$  \*/  
 $P_q \leftarrow P_q \cup \{ R_p \}$ ;  
 }  
 /\*  $A(B_p)$  is not identical with any  $A(B)$  \*/  
**else** {  
 $q++$ ;  
 $P_q = \{ R_p \}$ ;  
 $\mathcal{P} \leftarrow \{ P_q \}$ ;  
 }  
 }  
 }  
 }  
 /\* Generate new attribute candidate information  $AN_q$  from each  $P_q$  \*/  
**for**( $q = 1; q \leq |\mathcal{P}|; q++$ ) {  
 $AN_q = \langle \mathcal{A}_q, \mathcal{V}_q, \mathcal{S}_q, \mathcal{C}_q \rangle$ ;  
 $\mathcal{A}_q = A(B)$ ;  
 $\mathcal{V}_q = \{ V_{q_0}, V_{q_1}, \dots, V_{q_{|P_q|}} \}$ ;  
 where  
 $V_{q_0} = \bigwedge_{\forall R_p \in P_q} \neg true(B_k)$ ;  
 $V_{q_k} \in \{ true(B_p) | \forall R_p \in P_q \}$  ( $k = 1, \dots, |P_q|$ );  
 $\mathcal{S}_q = \{ sup(R_p) | \forall R_p \in P_q \}$ ;  
 $\mathcal{C}_q = \{ conf(R_p) | \forall R_p \in P_q \}$ ;  
 $\mathcal{AN} \leftarrow \mathcal{AN} \cup \{ AN_q \}$ ;  
 }  
 }  
 }

---

merged into the new attribute of  $\mathcal{A}_q$ . The value  $\mathcal{V}_q$  of a new attribute candidate is defined by the predicate ‘*true*( $B$ )’. This predicate becomes true when all of the items in Body  $B$  appear in an instance. The set of support and confidence values of AARs  $R \in P_q$  that constitutes the new attribute candidates are recorded as  $\mathcal{S}_q$  and  $\mathcal{C}_q$  for the evaluation of the candidate in a later step.

The proposed attribute generation process is explained using the following

example. The original training data consists of two attributes,  $A_1$  (attribute value  $V_1 = \{0, 1\}$ ),  $A_2$  (attribute value  $V_2 = \{0, 1\}$ ), and a class  $C$  (class  $C = \{0, 1\}$ ). Furthermore, suppose that the following two AARs are extracted:

$$R_1 : \mathbf{if} \{ \langle A_1, 0 \rangle, \langle A_2, 0 \rangle \} \mathbf{then} \{ \langle C, 0 \rangle \},$$

$$\mathit{support} : \mathit{sup}(R_1), \mathit{confidence} : \mathit{conf}(R_1), \quad (8)$$

where

$$\mathbf{B}_1 = \{ \langle A_1, 0 \rangle, \langle A_2, 0 \rangle \},$$

$$\mathbf{H}_1 = \{ \langle C, 0 \rangle \}, \text{ and}$$

$$A(\mathbf{B}_1) = \{ A_1, A_2 \}.$$

$$R_2 : \mathbf{if} \{ \langle A_1, 1 \rangle, \langle A_2, 1 \rangle \} \mathbf{then} \{ \langle C, 0 \rangle \},$$

$$\mathit{support} : \mathit{sup}(R_2), \mathit{confidence} : \mathit{conf}(R_2), \quad (9)$$

where

$$\mathbf{B}_2 = \{ \langle A_1, 1 \rangle, \langle A_2, 1 \rangle \},$$

$$\mathbf{H}_2 = \{ \langle C, 0 \rangle \}, \text{ and}$$

$$A(\mathbf{B}_2) = \{ A_1, A_2 \}.$$

The two AARs  $R_1$  and  $R_2$  have identical  $A(\mathbf{B})$ , i.e.,  $A(\mathbf{B}_1)$  and  $A(\mathbf{B}_2)$ . Therefore, these AARs are in the same partition element  $P_1 = \{R_1, R_2\}$ . A new attribute candidate  $AN_1$  is generated as follows:

$$AN_1 = \langle \mathcal{A}_1, \mathcal{V}_1, \mathcal{S}_1, \mathcal{C}_1 \rangle, \quad (10)$$

where

$$\mathcal{A}_1 = \{ A_1, A_2 \},$$

$$\mathcal{V}_1 = \{ V_{1_0}, V_{1_1}, V_{1_2} \},$$

$$\mathcal{S}_1 = \{ \mathit{sup}(R_1), \mathit{sup}(R_2) \}, \text{ and}$$

$$\mathcal{C}_1 = \{ \mathit{conf}(R_1), \mathit{sup}(R_2) \}.$$

Here,  $V_{1_1} = \mathit{true}(\mathbf{B}_1)$ ,  $V_{1_2} = \mathit{true}(\mathbf{B}_2)$ , and  $V_{1_0} = \bigwedge_{\forall R \in P_1} \neg \mathit{true}(\mathbf{B})$ . The attribute generation process explained above is depicted in Fig. 1.

### 3.4. Step 4: Evaluation of Generated Candidates

To evaluate the goodness of generated candidates to induce a decision tree, we adopt the information gain criterion used in the decision tree algorithm ID3 (Russell and Norvig, 1995). We use an approximated definition of gain,  $\mathit{Gain}(\mathcal{A}_q)$ , to reduce the computation cost for the large amount of data. In the evaluation of the approximated gain, the information given in the former steps, e.g., support, confidence, and total number of instances, is used, but any other information requiring further heavy computation is not needed. This feature of the approximation enhances the applicability of our proposed method to the large amount of data through a significant reduction in computation time.

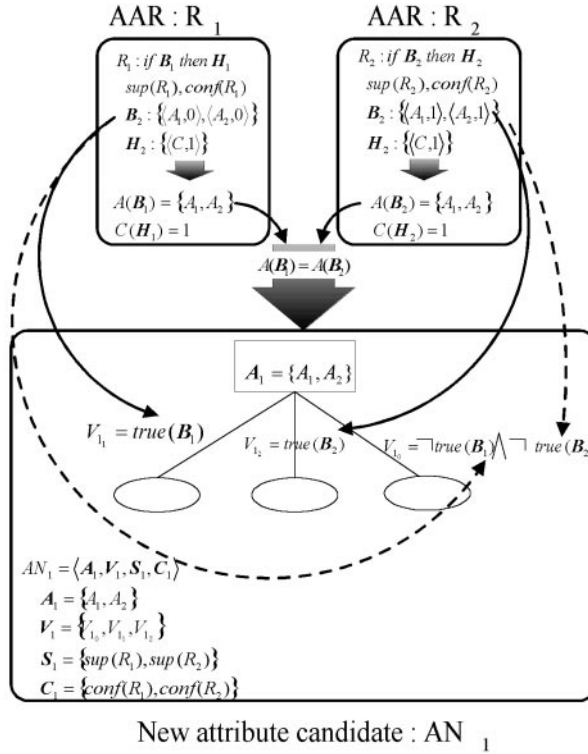


Fig. 1. New attribute candidate generation from AARs : R<sub>1</sub> and R<sub>2</sub>.

In detail, the information gain,  $Gain(\mathcal{A}_q)$ , is calculated by using the support values in  $\mathcal{S}_q$  and the confidence values in  $\mathcal{C}_q$ . Let  $R_k \in P_q$  ( $k = 1, \dots, |P_q|$ ).

$$\begin{aligned}
 & Gain(\mathcal{A}_q) \\
 &= - \sum_{j=1}^{|n_c|} \frac{n(v_{c,j})}{N} \log_2 \frac{n(v_{c,j})}{N} \\
 & - \sum_{k=1}^{|P_q|} \frac{sup(R_k)}{conf(R_k)} \left\{ -conf(R_k) \log_2 conf(R_k) - p_k \log_2 \frac{p_k}{n_c - 1} \right\} \\
 & - \left\{ 1.0 - \sum_{k=1}^{|P_q|} \frac{sup(R_k)}{conf(R_k)} \right\} \{-r_j \log_2 r_j\}, \tag{11}
 \end{aligned}$$

where

- $N$  : number of training data,
- $sup(R_k) \in \mathcal{S}_q$  : support value of  $R_k$  which corresponds to  $V_{qk}$ ,
- $conf(R_k) \in \mathcal{C}_q$  : confidence value of  $R_k$  which corresponds to  $V_{qk}$ ,
- $p_k = 1.0 - conf(R_k)$ ,



$$r_j = \frac{N_{q_0}(v_{c,j})}{N - N \sum_{k=1}^{|P_q|} \frac{sup(R_k)}{conf(R_k)}},$$

$n(v_{c,j})$  : number of instances whose class is  $v_{c,j}$ ,

$$N_{q_0}(v_{c,j}) = n(v_{c,j}) - \sum_{k=1}^{|P_q|} N_{q_k}(v_{c,j}),$$

: number of instances whose attribute value and class value are  $V_{q_0}$  and  $v_{c,j}$  respectively,

$$N_{q_k}(v_{c,j}) = \begin{cases} N \cdot sup(R_k) & (v_{c,j} \equiv C(H_k)) \\ \frac{N \cdot \frac{sup(R_k)}{conf(R_k)} (1.0 - conf(R_k))}{n_c - 1} & (\text{otherwise}) \end{cases}.$$

In this equation, the first term of the right-hand side represents the amount of information necessary to predict the class of the instances without using any information of the attributes. The second term represents the sum of the information amount needed to predict the classes of the instances having the attribute value  $V_{q_k} \in \mathcal{V}_q$ ,  $k = 1, \dots, |P_q|$ . The third term stands for the amount of information needed to classify the instances having the attribute value  $V_{q_0}$ . For more detail, refer to the Appendix.

If this information gain of the candidate is larger than 0, the candidate is considered to be informative for classification.

### 3.5. Step 5: Adding New Attributes

The new attribute candidates  $\mathcal{A}$  whose information gains  $Gain(\mathcal{A}_q)$  are larger than 0 are added to the attributes of the original data. By adding these new attributes to the data, the descriptive power of the data is expected to be improved, and a better decision tree may be induced.

## 4. Experiment

### 4.1. Conditions of Experiment

To confirm the effectiveness of our proposed method for the improvement of decision tree induction, experiments have been conducted for several sample data sets. In these experiments, we use C4.5 for decision tree induction (Quinlan, 1993). All of the functional options of C4.5 are set to defaults. The pruned decision trees are used for evaluation.

The sample data sets are selected from the UCI Machine Learning Repository (Blake et al., 1998). They are selected from data sets used in Zheng (1995) to compare the results. The specifications of the test data sets are summarized in Table 2. Because the proposed method can only deal with nominal attributes, only the data sets that include nominal attributes are used. The Monk's data sets (Monk1, Monk2, Monk3) comprise artificial data designed for evaluation of machine-learning algorithms. The other five data sets are real-world domain data. They are from a molecular biology domain (Promoters), three linguistic domains (Phoneme, Stress, Letter), and a game domain (Tic-Tac-Toe).

**Table 2.** The specifications of data sets for experiment

Data set	# of Training data	# of Test data	# of Attribute	# of Class
Monk1	124	432	6	2
Monk2	169	432	6	2
Monk3	122	432	6	2
Promoters	8124	CV10	22	2
Phenome	12960	CV10	8	52
Stress	12960	CV10	8	5
Letter	12960	CV10	8	163
Tic-Tac-Toe	958	CV10	9	2

\*CV10, 10-fold cross-validation.

**Table 3.** Accuracy (%) and size on the Monk's data sets

Algorithm	Monk1		Monk2		Monk3	
	Accuracy	Size	Accuracy	Size	Accuracy	Size
C4.5	75.7	18	65.0	31	97.2	12
AQ17-DCI	100.0	N/A	100.0	N/A	94.2	N/A
AQ17-HCI	100.0	N/A	93.1	N/A	100.0	N/A
CI	100.0	14	67.1	22	95.8	14
ID2-of-3	100.0	18	98.1	24	97.2	21
XofN	100.0	17	100.0	13	100.0	9
AARs	100.0	8	75.5	35	92.8	11

We evaluate mainly the following two aspects for data mining:

- *Improvement of decision tree*: Two indices are used to evaluate the effect of the proposed method for improvement of the decision tree. One is the size of the decision tree, and the other is the prediction accuracy.
- *Applicability to large data sets*: The applicability to the large data is evaluated. In the experiment, we investigate the computation time for various data sizes.

The features of the decision tree are evaluated using 10-fold cross-validation, except for the Monk's data sets.

## 4.2. Effect to Improve Decision Tree

The experimental results on Monk's data sets of our proposed method (AARs method) are summarized in Table 3. For comparison, we also give the results of some other attribute generation methods: AQ17-DCI, AQ17-HCI (Bloedorn and Michalski, 1998), CI (Zheng, 1992), ID2-of-3 (Murphy, 1991), XofN (Zheng, 1995) and decision trees generated with original attribute (C4.5). In this experiment, we set the minimum support and minimum confidence at 0.05 and 0.95 respectively.

In the Monk1 data set, both the size and the prediction accuracy are much improved by the AARs method. On the other hand, the AARs method does not demonstrate good performance in the Monk2 and Monk3 data sets.

Next, we evaluate the experimental results using real-world domain data. The experimental results are summarized in Tables 5 and 6. For the AARs method, we also investigate performance with several settings of minimum confidence and minimum support. The result of the AARs method is chosen for the case where

**Table 4.** Settings of minimum support and minimum confidence when best accuracy is demonstrated

Data set	AARs	
	Minimum support	Minimum confidence
Promoters	0.25	0.95
Phoneme	0.001	0.80
Stress	0.01	0.85
Letter	0.003	0.90
Tic-Tac-Toe	0.05	0.45

**Table 5.** Accuracy (%) on real-world domains

Data set	C4.5	CI	ID2-of-3	XofN	AARs
Promoters	76.3	81.0	87.6	88.5	90.4
Phoneme	81.1	82.3	83.1	83.9	83.6
Stress	82.7	83.8	86.2	87.6	86.1
Letter	73.7	65.6	75.1	76.9	76.6
Tic-Tac-Toe	84.7	94.2	94.9	98.4	99.7

the induced decision tree's prediction accuracy is the best among several settings. The best setting of minimum support and minimum confidence is summarized in Table 4.

In Promoter and Tic-Tac-Toe data sets, the AARs method improves both prediction accuracy and size over C4.5. On the other hand, in the three linguistic data sets (Phoneme, Letter, Stress), the AARs method cannot improve the size although prediction accuracy is improved compared to C4.5.

### 4.3. Applicability to Large-Scale Data

Next, the applicability of the proposed method to large-scale data is evaluated. One of the main factors that affect the computation time of the proposed method is the number of training data. Various sizes of test data sets have been generated as follows. An instance is picked up from Monk1 data set, and the class is changed to an erroneous value with probability 5%. This process is repeated until the number of data values needed is attained. Training data sets of sizes 10,000, 50,000, 100,000, and 500,000 are prepared. A personal computer having the specification of Linux OS, Pentium 166 Hz CPU and main memory 128 Mbytes is used in this experiment. The experimental results are shown in Table 7, from which the following are concluded:

**Table 6.** Size on the real-world domains

Data set	C4.5	CI	ID2-of-3	XofN	AARs
Promoters	22.6	15.0	11.2	13.9	7.0
Phoneme	2339.2	1634.5	1188.4	1506.0	3221.6
Stress	2077.3	1074.1	961.5	739.6	1751.6
Letter	3394.9	1024.9	1654.8	2242.4	2825.8
Tic-Tac-Toe	128.5	82.0	95.8	42.8	23.7

**Table 7.** The number of training data values and computation time for pre-processing

# of Training Data	C4.5 (Original data)		AARs (Pre-processed data)		Pre-processing computation time (s)
	Accuracy (%)	Size	Accuracy (%)	Size	
	10,000	78.7	90	95.0	
50,000	82.9	79	95.1	12	13
100,000	85.4	79	95.0	12	22
500,000	80.3	90	95.0	13	114

- The computation time for data pre-processing also increases when the data size becomes large. The increase is almost proportional to the size of the data.
- The effect of improving the decision tree is maintained even if the data size is large.

## 5. Discussion

### 5.1. Basic Features of the Proposed Method

The basic features of the proposed method are well demonstrated in the case of the decision tree induced from Monk's data sets. Monk's data sets are artificial data prepared for classification problem (Thrun et al., 1991). As indicated in the experimental results in Section 4.2, the AARs method improves the prediction accuracy and size of the induced decision tree over the other attribute generation methods. However, the AARs method does not demonstrate good performance in Monk2 and Monk3 data sets. We explain the reason for the performance difference.

Each data set has six attributes,  $a_1, \dots, a_6$ , and two class values,  $Class = \{0, 1\}$ . Each data set contains its target concept. The logical descriptions of the target concept are as follows:

**Monk1:** If  $(a_1 = a_2)$  or  $(a_5 = 1)$  then  $Class = 1$ .

**Monk2:** If  $(a_n = firstvalue)$  for exactly two choices of  $n$  in  $\{1, 2, \dots, 6\}$  then  $Class = 1$ .

**Monk3:** If either  $(a_5 = 3$  and  $a_4 = 1)$  or  $(a_5 \neq 3$  and  $a_1 \neq 3)$  then  $Class = 1$ . Monk3 data has 5% additional noise (misclassifications) in the training data.

Here, ' $a_x = a_y$ ' means that these attributes take the same attribute value. The performance of the proposed method is summarized for each data as follows:

- *Monk1:* By applying the AARs method to the original training data, 16 new attributes  $a_7, \dots, a_{22}$  are generated. The decision trees induced from the original data and the pre-processed data by the AARs method are depicted in Fig. 2 respectively. The decision tree induced from the training data pre-processed by the AARs method includes new attributes  $a_7$  in the root node, and  $a_{10}$  as a node at the next level. The new attribute  $a_7$  implies a target concept 'if  $(a_1 = a_2)$  then  $Class = 1$ ' appropriately with conjunction of two original attributes. These included new attributes imply the target concept, so that the AARs method improves the prediction accuracy and size of decision trees.

- *Monk2*: Thirteen new attributes  $a_7, \dots, a_{19}$  are generated. In this example, we notice the limitation of the AARs method and its new attribute representation. As new attributes, the AARs method can generate only conjunction (not disjunction) of original attributes. Even under this representation constraint on new attributes, the proposed method can extract the correct target concept of the data in the form of AARs such as

**if**  $\{\langle a_1, 1 \rangle, \langle a_2, 1 \rangle, \langle a_3, 0 \rangle, \langle a_4, 0 \rangle, \langle a_5, 0 \rangle, \langle a_6, 0 \rangle\}$  **then**  $\{Class = 1\}$

...

**if**  $\{\langle a_1, 0 \rangle, \langle a_2, 0 \rangle, \langle a_3, 0 \rangle, \langle a_4, 0 \rangle, \langle a_5, 1 \rangle, \langle a_6, 1 \rangle\}$  **then**  $\{Class = 1\}$

However, such conjunctive associations of many attributes rarely appear; i.e., their support is very small in this data set. For this reason, extracting such associations and generating new attributes are very difficult for this data set. Like the AARs method, CI (Zheng, 1992) also represents new attributes as a conjunction of original attributes. Therefore, CI does not improve the decision tree in Monk2.

On the other hand, the XofN and its new attribute representation X-of-N (Zheng, 1995) can represent the target concept as follows:

**if**  $X - of - \{\langle a_1, 1 \rangle, \langle a_2, 1 \rangle, \langle a_3, 0 \rangle, \langle a_4, 1 \rangle, \langle a_5, 1 \rangle, \langle a_6, 1 \rangle\} = 2$   
**then**  $\{Class = 1\}$

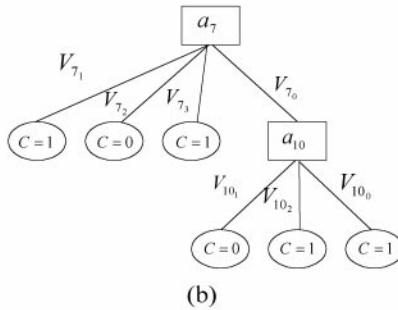
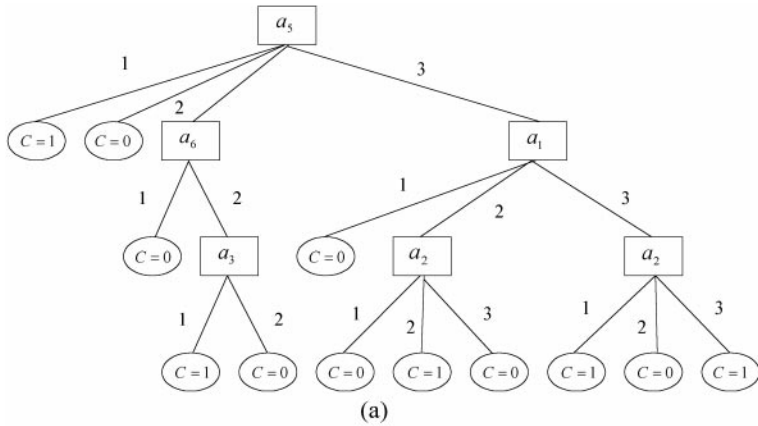
The target concept of Monk2 is more suitable for X-of-N representation than AARs' representation. This is one of the major reasons why the improvement of the decision tree is different between XofN and the AARs method.

- *Monk3*: The 21 new attributes  $a_7, \dots, a_{27}$  are generated. The decision tree from the data pre-processed by the AARs method includes new attributes  $a_{13}$ . By data pre-processing, the size of the decision tree becomes smaller. However, the prediction error is slightly increased. The proposed method extracts one of the target concepts '**if** ( $a_5 = 3$  and  $a_4 = 1$ ) **then**  $Class = 1$ '. However, the other target concept '**if** ( $a_5 \neq 3$  and  $a_1 \neq 3$ ) **then**  $Class = 1$ ' is not extracted by AARs because of the representation constraint mentioned above. Owing to this limitation of extracting the latter target concept, the good set of new attributes which expresses the target concept appropriately cannot be generated. This fact adversely affects prediction accuracy of the decision tree induced from the data pre-processed by the AARs method.

## 5.2. Effectiveness of the Proposed Method in Data Mining

In most of the data sets, the AARs method improves prediction accuracy. The AARs method improves both prediction accuracy and size particularly in Monk1, Promoter, and Tic-Tac-Toe. The target concepts in these data sets are suitable for representation using AARs. On the other hand, the AARs method does not work well in Monk2 and Monk3 data sets, where target concepts are not suitable for the AARs method's new attribute representation. These target concepts are more suitable for representation using X-of-N and negation.

The adopted new attribute representation is different from other attribute generation methods. From the practical point of view, attribute generation meth-



$a_7, a_{10}$  : new attributes

$$a_7 = \{a_1, a_2\}$$

$$V_{71} = \text{true}(\{\{a_1,1\}, \{a_2,1\}\})$$

$$V_{72} = \text{true}(\{\{a_1,2\}, \{a_2,2\}\})$$

$$V_{73} = \text{true}(\{\{a_1,3\}, \{a_2,3\}\})$$

$$V_{70} = \bigwedge_{k=1}^3 \neg V_{7k}$$

$$a_{10} = \{a_3, a_2\}$$

$$V_{101} = \text{true}(\{\{a_3,1\}, \{a_2,1\}\})$$

$$V_{102} = \text{true}(\{\{a_3,2\}, \{a_2,2\}\})$$

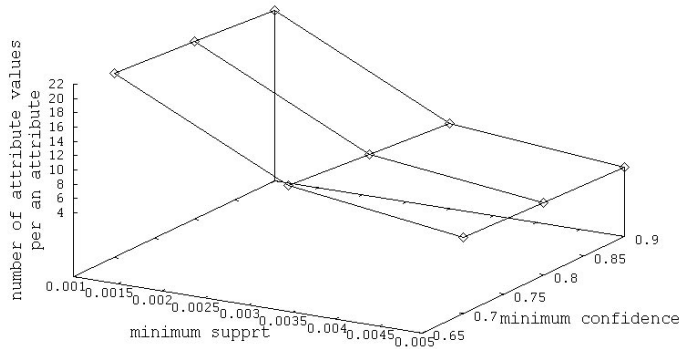
$$V_{100} = \bigwedge_{k=1}^2 \neg V_{10k}$$

**Fig. 2.** Two decision trees induced from the Monk1 data set: (a) original; (b) pre-processed by the AARs method.

ods which adopt different representations should be applied to the original data, and the generated attributes should be selected by an attribute selection method.

When the training data size becomes large, the computation time required by the proposed method increases. However, the increase is almost proportional to the size of the data. This result guarantees the applicability of the proposed method to large data sets.

The levels of minimum support and minimum confidence are main parameters of the proposed method for new attribute generation. When the minimum support is set small, the associations among attributes which appear infrequently in the data are extracted as AARs. On the other hand, association rules which are not plausible are extracted as AARs when the minimum confidence is set small. Thus, the smaller the values of the minimum support and the minimum confidence, the



**Fig. 3.** The relation between minimum thresholds and the number of attribute values per new attributes (Phenome).

more AARs are extracted. Accordingly, each new attribute takes many attribute values. The relations between these parameters' thresholds and the number of attribute values per new attribute in the Phoneme data set is depicted in Fig. 3. By including such new attributes taking many attribute values, the size of the decision tree becomes larger.

The association pattern appears less frequently when either the number of attribute values or the number of classes increases. Therefore, the minimum support should be set smaller when the number of attribute values per attribute or number of classes is larger. Figure 4 depicts the relationship between (a) the product of the number of attribute values per attribute and the number of classes and (b) the minimum support value when the prediction accuracy of the induced decision tree is the best. The results are for the real-world data used in Section 4.2. The negative correlation is confirmed between these two values. The regression equation is

$$y = 0.34x^{-0.67} \quad (R^2 = 0.820), \tag{12}$$

where

$x$  : product of the number of attribute values per attribute and the number of classes,

$y$  : minimum confidence when the prediction accuracy of the induced decision tree is the best, and

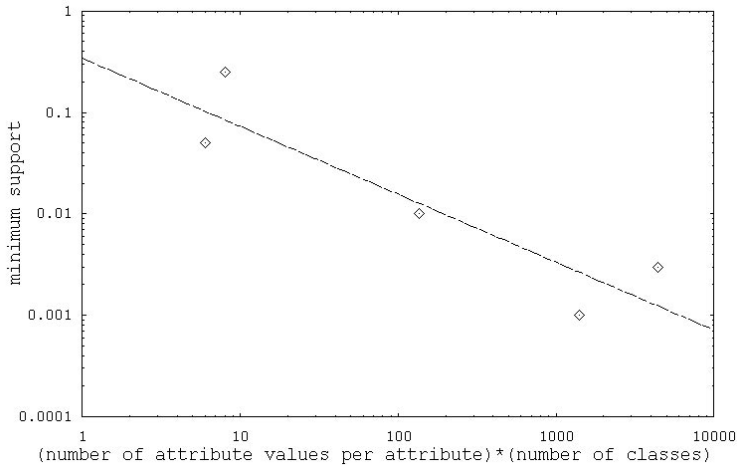
$R^2$  : coefficient of determination

The coefficient of determination  $R^2$  is calculated as follows. This represents the fraction of the deviation component explained by this regression in the data.

$$R^2 = \frac{\sum_{i=1}^n (\hat{y}_i - \bar{y})^2}{\sum_{i=1}^n (y_i - \bar{y})^2}, \tag{13}$$

where

$\{x_i, y_i\}$  : (observed) data ( $i = 1, \dots, n$ ),



**Fig. 4.** The relationship between  $x$  (product of the number of attribute values per attribute and classes) and  $y$  (the minimum support value when the prediction accuracy of induced decision tree is the best).

$\hat{y}_i$  : estimates on  $y_i$  predicted with the regression equation, and  
 $\bar{y}$  : mean of  $y$

The regression equation indicates that the minimum support should be set smaller when the product of the number of attribute values and classes is larger.

## 6. Related Work

Many have investigated attribute generation methods. ID2-of-3 (Murphy and Pazzani, 1991) adopts M-of-N new attribute representation, and generates binary attributes. XofN (Zheng, 1995), an extension of ID2-of-3, represents new nominal attributes by X-of-N. As discussed in Section 5, these methods improve the prediction accuracy and size when these new attribute representations are appropriate to describe the target concept. As for representation of new attributes, CI (Zheng, 1992) is one of the closest to the AARs method, because it also uses conjunction to represent the target concept of data. However, the authors did not discuss the applicability to data mining dealing with large-volume data.

Bloedorn and Michalski (1998) investigated the method of data-driven constructive induction (DCI). In this method, new attributes are generated by applying the various attribute construction operators to the original ones. Although evaluation of the new attributes is conducted by using training data, the operators need to be selected from the prepared repository by the human user.

Lavráč et al. (1998) investigated the method of attribute generation based on a priori knowledge of attributes in the framework of inductive logic programming. Although the algorithm works very well in many cases, applicability to large-scale data has not been assessed.

Liu et al. (1998) proposed a classifier induction algorithm which induces the classifiers from association rules. The algorithm shows better performance on



the classification problem than C4.5. However, their algorithm does not induce a decision tree but only a rule classifier.

## 7. Conclusion

In this paper, we proposed a novel data pre-processing method to improve the performance of decision tree induction. The effectiveness of our method has been demonstrated through experiments using a subset of UC Irvine data sets. Without using a priori knowledge of attributes and associations among them, our method extracts the knowledge of association among attributes in the form of AARs from the data automatically, and uses them to generate new attributes. The computation time required by the method remains small even if the size of the training data is large. These features are highly advantageous for the purpose of data mining.

The following issues remain for our future work:

- By applying the proposed method to the decision tree algorithm such as C4.5, we expect to induce appropriate decision trees. We plan to evaluate the effect to improve decision trees by our method together with the attribute grouping function implemented by C4.5 (Quinlan, 1993).
- Currently, our method is not directly applicable to continuous valued attributes. Adopting a function dealing with continuous valued attributes to our method will be a future study.
- In processing new attribute evaluation by equation (11), we mainly consider processing time. The process does not need much calculation cost, and this feature is appropriate for the data-mining method. However, attribute selection mechanisms should be investigated because increase of the number of attributes requires the expansion of the database and more calculation to induce the decision tree.

## References

- Agrawal R, Srikant R. (1994) Fast algorithms for mining association rules. In Proceedings of the 20th VLDB Conference, Morgan Kaufmann, San Mateo, CA, pp 487–499
- Berry MJA, Linoff GS (1997) Data mining techniques for marketing, sales, and customer support. Wiley, New York
- Blake C, Keogh E, Merz CJ (1998) UCI repository of machine learning databases, <http://www.ics.uci.edu/mllearn/MLRepository.html>, University of California, Department of Information and Computer Science, Irvine, CA
- Bloedorn E, Wnek J, Michalski RS. (1993) Multistrategy constructive induction AQ17-MCI. In Proceedings of the 2nd international workshop on multistrategy learning, Morgan Kaufmann, San Mateo, CA, pp 188–203
- Bloedorn E, Michalski RS. (1998) Data-driven constructive induction: a method and its application. IEEE Intelligent Systems and their Applications 13(2): 30–37
- Lavrác N, Gamberger D, Turney P (1998) A relevancy filter for constructive induction. IEEE Intelligent Systems and their Applications 13(2): 50–56
- Liu H, Motoda H (eds) (1998) Feature selection for knowledge discovery and data mining. Kluwer, Dordrecht
- Liu B, Hsu W, Ma Y (1998) Integrating classification and association rule mining. In Proceedings of the 4th conference on knowledge discovery and data mining, AAAI Press, Menlo Park, CA, pp 80–86
- Murphy PM, Pazzani MJ (1991) ID2-of-3: constructive induction of M-of-N concepts for discriminators in decision trees. In Proceedings of the 8th international workshop on machine learning, Morgan Kaufmann, San Mateo, CA, pp 183–187

- Quinlan R (1993) C4.5: programs for machine learning. Morgan Kaufmann, San Mateo, CA
- Russell S, Norvig P (1995) Artificial intelligence: a modern approach. Prentice-Hall, Englewood Cliffs, NJ
- Thrun S, Bala J, Bloedorn E et al. (1991) The MONK's problems: a performance comparison of different learning algorithms. Technical Report of Carnegie Mellon University, CMU-CS-91-197
- Zheng Z (1992) Constructing conjunctive tests for decision trees. In Proceedings of the 5th Australian joint conference on artificial intelligence, World Scientific, Singapore, pp 355–360
- Zheng Z (1995) Constructing nominal X-of-N attributes. In Proceedings of the 14th international joint conference of artificial intelligence, Morgan Kaufmann, San Mateo, CA, pp 1064–1070

## Appendix: On Evaluation Equation (11)

We give additional explanation on how to derive the evaluation equation (11) for a new attribute candidate  $\mathcal{A}_q$ . Here, its attribute values  $\mathcal{V}_q$  are assumed to be  $\mathcal{V}_q = \{V_{q_1}, \dots, V_{q_{|P_q|}}, V_{q_0}\}$ .

The new attribute candidate is evaluated using the approximated information gain  $Gain(\mathcal{A}_j)$  when it is adopted at the root node of the decision tree.

First, the amount of information  $G_{root}$ , which is necessary to predict the class correctly in the root node, is calculated with the number of instances  $n(v_{c,j})$  which belong to each class  $v_{c,j}$ :

$$G_{root} = - \sum_{j=1}^{n_c} \frac{n(v_{c,j})}{N} \log_2 \frac{n(v_{c,j})}{N}, \quad (14)$$

where

$N$  : number of training data,

$n(v_{c,j})$  : number of instances whose class is  $v_{c,j}$

Next, the amount of information which is necessary to predict the class in the child node  $child_k$ , which has the following the attribute value  $\{V_{q_k} | k = 1, \dots, |P_q|\} \in \mathcal{V}_q$  defined by the AAR  $R_k$ , is calculated as follows.

From the definition of support value  $sup(R_k)$  and confidence value  $conf(R_k)$  of  $R_k$ , the number  $N_{q_k}^1$  of instances falling into  $child_k$  is equal to the number of instances which match the condition part of AAR  $R_k$ . Therefore,  $N_{q_k}^1$  is calculated with support value  $sup(R_k)$  and confidence value  $conf(R_k)$  as follows:

$$N_{q_k}^1 = N \cdot \frac{sup(R_k)}{conf(R_k)} \quad (15)$$

Furthermore, the number of instance  $N_{q_k}^2$  which belongs to the class  $C(H_k)$  is calculated by the support value  $sup(R_k)$  of  $R_k$ :

$$N_{q_k}^2 = N \cdot sup(R_k) \quad (16)$$

On the other hand, the number of instances which belong to other classes cannot be calculated using the index of AARs alone. Thus, we take the conservative approximation by considering the case that the amount of information to predict the class correctly is the maximum. The case occurs when the numbers of instances that belong to each class are even. The number of instance  $N_{q_k}^3$  is

calculated as follows:

$$\begin{aligned} N_{q_k}^3 &= \frac{N_{q_k}^1 - N_{q_k}^2}{n_c - 1} \\ &= \frac{N \cdot \frac{\text{sup}(R_k)}{\text{conf}(R_k)} (1.0 - \text{conf}(R_k))}{n_c - 1} \end{aligned} \quad (17)$$

Hence, the approximated amount of information  $G_{child_k}$  which is necessary to predict the class correctly in the child node  $child_k$  can be estimated as follows:

$$\begin{aligned} G_{child_k} &= \frac{N_{q_k}^1}{N} \left\{ - \left( \frac{N_{q_k}^2}{N_{q_k}^1} \right) \log_2 \left( \frac{N_{q_k}^2}{N_{q_k}^1} \right) - (n_c - 1) \cdot \left( \frac{N_{q_k}^3}{N_{q_k}^1} \right) \log_2 \left( \frac{N_{q_k}^3}{N_{q_k}^1} \right) \right\} \\ &= \frac{\text{sup}(R_k)}{\text{conf}(R_k)} \left\{ -\text{conf}(R_k) \log_2 \text{conf}(R_k) - p_k \log_2 \frac{p_k}{n_c - 1} \right\}, \end{aligned} \quad (18)$$

where

$$p_k = 1.0 - \text{conf}(R_k)$$

$N_{q_0}^1$ , the number of instances which belong to the child node  $child_0$  corresponding to the attribute value  $V_{q_0}$ , is calculated by subtracting the number of instance belonging to the other child node  $child_k$  from the number of all instance  $N$ :

$$\begin{aligned} N_{q_0}^1 &= N - \sum_{k=1}^{|P_q|} N_{q_k}^1 \\ &= N \left\{ 1.0 - \sum_{k=1}^{|P_q|} \frac{\text{sup}(R_k)}{\text{conf}(R_k)} \right\} \end{aligned} \quad (19)$$

$N_{q_0}(v_{c,j})$ , the number of instances which belong to the class  $v_{c,j}$  is estimated by the number of instances of the same class in other children:

$$N_{q_0}(v_{c,j}) = n(v_{c,j}) - N_q(v_{c,j}), \quad (20)$$

where

$N_q(v_{c,j})$  : number of instances whose attribute values are  $\{V_{q_k} | k = 1, \dots, |P_q|\}$ , and class is  $v_{c,j}$ ,

$$N_q(v_{c,j}) = \sum_{k=1}^{|P_q|} N_{q_k}(v_{c,j}),$$

$N_{q_k}(v_{c,j})$  : number of instances whose attribute values are  $V_{q_k}$  and class is  $v_{c,j}$ ,

$$N_{q_k}(v_{c,j}) = \begin{cases} N \cdot \frac{\text{sup}(R_k)}{\text{conf}(R_k)} & (v_{c,j} \equiv C(\mathbf{H}_k)) \\ \frac{N \cdot \frac{\text{sup}(R_k)}{\text{conf}(R_k)} (1.0 - \text{conf}(R_k))}{n_c - 1} & (\text{otherwise}) \end{cases}$$

As mentioned above, the information quantity  $G_{child_0}$  which is necessary to predict

the class correctly in the child node  $child_0$  is calculated as follows:

$$\begin{aligned}
 G_{child_0} &= \frac{N_{q_0}^1}{N} \sum_{j=1}^{n_c} \left\{ - \left( \frac{N_{q_0}(v_{c,j})}{N_{q_0}^1} \right) \log_2 \left( \frac{N_{q_0}(v_{c,j})}{N_{q_0}^1} \right) \right\} \\
 &= \left\{ 1.0 - \sum_{k=1}^{|P_q|} \frac{sup(R_k)}{conf(R_k)} \right\} \sum_{j=1}^{n_c} \{-r_j \log_2 r_j\}, \\
 &\text{where} \\
 r_j &= \frac{N_{q_0}(v_{c,j})}{N - N \sum_{k=1}^{|P_q|} \frac{sup(R_k)}{conf(R_k)}}
 \end{aligned}$$

The approximated information gain  $Gain(\mathcal{A}_q)$  is therefore calculated by

$$Gain(\mathcal{A}_q) = G_{root} - \sum_{k=1}^{|P_q|} G_{child_k} - G_{child_0} \quad (21)$$

This leads to equation (11).

## Author Biographies



**Masahiro Terabe** is a staff researcher in the Department of Safety Science and Policy of Mitsubishi Research Institute, Inc., Japan. His current research interests include machine learning, knowledge acquisition, knowledge discovery and data mining. He received his B.S. (1993) and M.S. (1995) degrees in precision engineering from Kyoto University. He is a member of the AAAI, the Japanese Society for Artificial Intelligence, and the Society of Instrument and Control Engineers.



**Hiroshi Motoda** is a professor in the Division of Intelligent Systems Science at the Institute of Scientific and Industrial Research of Osaka University. Before joining the university, he had been with Hitachi since 1967 and reached the position of a senior chief research scientist at the Advanced Research Laboratory, where he headed an AI group and conducted research on machine learning, knowledge acquisition, diagrammatic reasoning and information filtering. His current research interests include, in addition to these, scientific knowledge discovery and data mining. He received his B.S. (1965), M.S. (1967) and Ph.D. (1972) degrees in nuclear engineering from the University of Tokyo. He is now on the editorial board of *Artificial Intelligence in Engineering*, *International Journal of Human-Computer Studies* and *Knowledge and Information Systems*. He received the outstanding achievement award from JSAI (1999). He is a member of the AAAI, IEEE Computer Society, JSAI, JSSST, IPSJ and JCSS.



**Takashi Washio** graduated at the Department of Nuclear Engineering, Tohoku University, in 1983, and took his M.S.E. and Ph.D. in the same department in 1985 and 1988 respectively. He was a visiting researcher in the Nuclear Reactor Laboratory of Massachusetts Institute of Technology from 1988 to 1990, and was a senior researcher in the Mitsubishi Research Institute, Inc. in Tokyo, Japan, from 1990 to 1996. He has researched on the techniques of qualitative reasoning, diagnosis theory and risk analysis for large-scale plants. He became an associate professor of the Institute for Scientific and Industrial Research (ISIR), Osaka University, in 1996. His current research interests are automated scientific law discovery and industrial data-mining techniques. He is the member of the AAAI, the Japanese Society for Artificial Intelligence, Society of Instrument and Control Engineers, Information Processing Society of Japan and Japan Society for Fuzzy Theory and Systems.



**Osamu Katai** graduated from Kyoto University at the Department of Mechanical Engineering in 1969, then proceeded to the Graduate School of Engineering and received M.S.E. and D.E. degrees in 1971 and 1979, respectively. Since 1971, he has been with Kyoto University, first as an Instructor at the Department of Precision Engineering and now is a Professor at the Department of Systems Science in the Graduate School of Informatics. From 1980 to 1981, he was a visiting researcher at INRIA (National Research Institute on Informatics and Automation, France). His current research interests are on the methodologies and theoretical frameworks of 'symbiotic systems' where artificial systems and natural systems can coexist harmoniously. Particular interests are on bio-informatic systems, co-evolving systems, autonomous robots, agents, environmental design, ecological design, community design, etc. He is a member of the Japanese Society for Artificial Intelligence and the Society of Instrument and Control Engineers, among others.



**Testuo Sawaragi** is an associate professor at the Department of Precision Engineering, Graduate School of Engineering, Kyoto University, Japan. He received his B.S., M.S. and Ph.D. degrees in Systems Engineering from Kyoto University in 1981, 1983 and 1988, respectively. From 1986 to 1994, he was an instructor at the Department of Precision Mechanics, Faculty of Engineering, Kyoto University, and in 1994 he was with the current department as an associate professor. From 1991 to 1992, he was a visiting scholar at the Department of Engineering-Economic Systems, Stanford University, USA. He has been engaged in research on Systems Engineering, Cognitive Science and Artificial Intelligence, particularly in the development of human-machine collaborative systems, modeling the transfer of human cognitive skills into machines. He is a member of the Society of Instrument and Control Engineers, the Institute of Systems, Control and Information Engineers, Japanese Society for Artificial Intelligence, Japan Society for Fuzzy Theory and Systems, Human Interface Society, JASME, and IEEE.

---

*Correspondence and offprint requests to:* Masahiro Terabe, Safety Science and Policy Department, Mitsubishi Research Institute, Inc., 2-6-1 Ohtemachi, Chiyoda, Tokyo 100-8141, Japan.  
Email:terabe@mri.co.jp