

Scalable Rule-Based Gene Expression Data Classification (Extended Version)

Mark A. Iwen
University of Michigan
Department of Mathematics
markiwen@umich.edu

Willis Lang
University of Michigan
EECS Department
wlang,jignesh@eecs.umich.edu

ABSTRACT

Current state-of-the-art association rule-based classifiers for gene expression data operate in two phases: (i) Association rule mining from training data followed by (ii) Classification of query data using the mined rules. In the worst case, these methods require an exponential search over the subset space of the training data set's samples and/or genes during at least one of these two phases. On large gene expression datasets these existing methods are computationally prohibitively expensive.

Our main result is the development of a heuristic rule-based gene expression data classifier called Boolean Structure Table Classification (BSTC). BSTC is explicitly related to association rule-based methods, but is guaranteed to be polynomial space/time. Extensive cross validation studies on several real gene expression datasets demonstrate that BSTC retains the classification accuracy of current association rule-based methods while being orders of magnitude faster than the current best classifier RCBT [9] on large datasets. As a result, *BSTC is able to finish table generation and classification on large datasets for which current association rule-based methods become computationally infeasible.*

BSTC also enjoys two other advantages over association rule-based classifiers: (i) BSTC is easy to use (requires no parameter tuning), and (ii) BSTC can easily handle datasets with any number of class types. Furthermore, in the process of developing BSTC we introduce a novel class of boolean association rules which have many potential applications to other data mining problems.

1. INTRODUCTION

Microarray technology allows biologists to simultaneously measure the expression of thousands of genes in a single experiment. This technology is currently revolutionizing biomedical research as it provides a unique tool to examine the state of a cell under various conditions. Microarray methods are also poised to play a critical role in personalized medicine as they can be used to determine the unique genetic susceptibility of an individual to disease.

See Table 1 for a sample microarray dataset shown using the common discretized relational representation. In this table, each sample row consists of (i) a list of discretized genes and (ii) a class

Sample	Expressed Genes				Class Label
s_1	g_1	g_2	g_3	g_5	Cancer
s_2	g_1	g_3	g_6		Cancer
s_3	g_2	g_4	g_6		Cancer
s_4	g_2	g_3	g_5		Healthy
s_5	g_3	g_4	g_5	g_6	Healthy

Table 1: Running Example of Microarray Data

label. A gene is present in a sample row if the sample expresses the gene. The absence of a gene in a row implies that the gene is not expressed in that sample. Thus, the sample/gene expression relationships for relational microarray data are essentially boolean.

Leading associated rule-based methods such as Top-k [9], FARMER [10], CLOSET+ [30], and CHARM [34] which have been applied to microarray datasets aim to correlate gene expression patterns with the classification labels. For these algorithms the discovered correlations take the form of **association rules** [2]. For an example association rule, consider the data shown in Table 1. Note that only the **Cancer** samples s_1 and s_2 express both genes g_1 and g_3 . Based on this observation we can create the following association rule: $g_1, g_3 \Rightarrow \text{Cancer}$. This rule implies that if a query sample express both g_1 and g_3 (i.e., if g_1 and g_3 's associated genes are both expressed in their associated expression intervals), then the query sample is likely to be of type **Cancer**. Hence, we can use this rule to *classify* query samples of unknown type as **Cancer** if they express both g_1 and g_3 . Note that there is nothing special about the class label **Cancer**. After noticing that only **Healthy** sample s_5 expresses both g_5 and g_6 , we can also create the meaningful association rule $g_5, g_6 \Rightarrow \text{Healthy}$.

In this paper we focus on association rule-based classifiers (hereafter referred to simply as rule-based classifiers) for gene expression data. We focus on rule-based classifiers for two reasons: (i) rule-based classifiers have been demonstrated to be more accurate for gene expression analysis than other methods [9, 10, 14, 18] such as SVM [11] and tree-based C4.5 family algorithms [26], and (ii) as opposed to other classifiers such as SVM, rule-based classifiers can offer concise, concrete, and biologically meaningful rules supporting their non-default classifications. However, rule-based methods are not scalable due to their high association rule mining costs. Although these rule mining costs are "one-time costs" in the sense that rules must only be mined once per training set, larger training data sets are being generated at an ever increasing rate. It is impossible for any exponential time method to keep up. Consequently, in this paper, we focus on extending accurate association rule-based classification methods to larger data sets.

This paper develops a scalable rule-based classifier called Boolean

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

ICDE '08 Cancún, México

Copyright 200X ACM X-XXXXX-XX-X/XX/XX ...\$5.00.

Structure Table Classification (BSTC) for microarray datasets. Given a labeled training set, such as the example in Table 1, BSTC *efficiently* builds an *accurate* classifier. The emphasis on accuracy is easy to appreciate and comes from BSTC being related to association rule-based methods. Hence, BSTC supports its classifications with intuitive rules. The emphasis on efficiency is also critical since large gene expression datasets are computationally taxing for existing association rule-based algorithms and, as successful microarray techniques fuel the growth of gene expression datasets, these methods will quickly become infeasible. In contrast, BSTC’s space and runtime costs are only polynomial. Hence, *BSTC is scalable to large data sets on which current association rule-based methods are computationally challenged.*

In an attempt to control runtime many current association rule methods [9, 10, 21] utilize support based rule pruning. Using a large enough support cutoff does allow rule mining to finish more quickly, but doesn’t completely resolve the issue. If the user sets the support cutoff too small he/she can easily spend days waiting for rule mining to finish before giving up in frustration. A few such mistakes can result in weeks of wasted time. On the other hand, setting the support cutoff too high excludes the generation of important high-confidence lower-support rules [22]. In order to not miss too many important rules the user can’t set the support cutoff too high. The end result is that in practice support cutoffs are difficult and time intensive to tune. In contrast, *BSTC is fast and easy to use.*

In addition, to the best of our knowledge all current association rule-based classifiers for gene expression data only handle datasets with two class labels. Although our example Table 1 data contains just two class labels, in practice microarray data can contain an arbitrary, though small, number of class types. *Unlike previous association rule-based classifiers, BSTC easily generalizes to datasets with more than two class types.*

To develop an accurate, scalable, multi-class, and easy to use rule-based classifier we carefully considered the underlying primitives that power association rule-based methods. These methods use **conjunctive association rules (CARs)**, where the rule antecedent is restricted to being a conjunction of terms. In contrast, we approach this problem by relaxing the types of rules to an important and larger subset of the more general class of boolean association rules (BARs). We develop a novel method for compactly storing these BARs in a simple data structure called a Boolean Structure Table (BST). BSTs can then be used for BAR generation and classification. BST classification (BSTC) collectively considers many simple BARs with 100% confidence in bulk. Because the rules are simple BSTC avoids extensive rule mining. Furthermore, considering rules in bulk keeps the computational cost low.

The main contributions of this paper are:

1. We propose a new polynomial time and space rule-based classifier for gene expression data analysis that is accurate, scalable, easy to use, and easily generalizable to multi-class classification.
2. We extensively evaluate our method against the current best association rule-based method (RCBT [9]), and show that our method is orders of magnitude faster on large datasets while maintaining high classification accuracy.
3. We introduce a subclass of more general boolean association rules and relate them to existing CARs. This not only leads to a better appreciation of why our classification method works, but also lays the foundation for the future use of these BARs on other database problems.

The remainder of this paper is organized as follows: First, in Sec-

tion 2 we formalize the concept of BARs. Then in Section 3, we define a concept called Boolean Structure Tables (BSTs) which are related to an important class of BARs. Section 5 provides a polynomial time and parameter-free classifier based on BSTs. Section 6 presents an extensive empirical evaluation of our classifier. Finally, Section 7 discusses related work, and Section 8 briefly presents our conclusions and directions for future work.

2. PRELIMINARIES

We work with the following type of data: We are given a finite set G of genes and N collections of subsets from G . These N collections are disjoint and represented as $C_1 = \{s_{1,1}, \dots, s_{1,m_1}\}, \dots, C_N = \{s_{N,1}, \dots, s_{N,m_N}\}$. Each C_i is called a **class type** or **class label**. Furthermore, we will refer to each set $s_{i,j} \subset G$ as a **sample** and every element $g \in G$ as a **gene**. We denote the total set of samples by $S = \cup_{i=1}^N C_i$. If $g \in s_{i,j}$ we will say that sample $s_{i,j}$ **expresses** gene g . Otherwise, if $g \in G$ and $g \notin s_{i,j}$ we will say that sample $s_{i,j}$ doesn’t express gene g . Similarly, we say that sample s is of class type C_i if and only if C_i contains $s \subset G$. Consider the Table 1 data. Here we have samples $S = \{s_1, s_2, s_3, s_4, s_5\}$ and genes $G = \{g_1, g_2, g_3, g_4, g_5, g_6\}$. Furthermore, we have $N = 2$ classes: $C_1 = \text{Cancer} = \{s_1, s_2, s_3\}$ and $C_2 = \text{Healthy} = \{s_4, s_5\}$.

Given such relational training data, a **conjunctive association rule (CAR)** is any element of $2^G \times \{1, \dots, N\}$. A CAR $g_{j_1}, \dots, g_{j_r} \Rightarrow n$ can be interpreted as follows: “If a query sample s contains all genes g_{j_1}, \dots, g_{j_r} , then it should be grouped with class type C_n .” Naturally, of the $2^{|G|} \cdot N$ possible association rules some are more useful than others. The following standard definitions were introduced in [2] to compare association rules:

Support: The support of a CAR $g_{j_1}, \dots, g_{j_r} \Rightarrow n$, called $supp[g_{j_1}, \dots, g_{j_r} \Rightarrow n]$, is:

$$|\{s_{n,j} \text{ s.t. } \{g_{j_1}, \dots, g_{j_r}\} \subset s_{n,j}, 1 \leq j \leq m_n\}|.$$

Confidence: The confidence of a CAR $(g_{j_1}, \dots, g_{j_r}, n)$ is:

$$\frac{supp[(g_{j_1}, \dots, g_{j_r}, n)]}{|\{s_{i,j} \text{ s.t. } \{g_{j_1}, \dots, g_{j_r}\} \subset s_{i,j} \forall i, j\}|}.$$

Consider the CAR $g_1, g_3 \Rightarrow \text{Cancer}$ for our running example in Table 1. We can see that the example CAR has a support of 2 since only two Cancer samples, s_1 and s_2 , contain both g_1 and g_3 . Furthermore, we can see that the example CAR has confidence 1 (or 100%) since no healthy samples contain both g_1 and g_3 .

2.1 Boolean Association Rules

For any sample s and gene g_i , $1 \leq i \leq n = |G|$, let $s[g_i] \in \{0, 1\}$ represent whether or not sample s expresses gene g_i . Furthermore, define $s[-g_i]$ to be $-s[g_i] \forall g_i \in G$. Now suppose that $B(x_1, \dots, x_n)$ is a Boolean expression whose value depends on some subset of $\{x_1, \dots, x_n\}$. We can evaluate B to true or false given a sample s by computing $B(s[g_1], \dots, s[g_n])$. For example, consider the boolean expression:

$$\hat{B}(x_1, x_2, x_3, x_4, x_5, x_6) = (x_1 \wedge x_3) \vee (x_2 \wedge x_4).$$

Using Table 1 we can evaluate

$$\hat{B}(s_1[g_1], s_1[g_2], s_1[g_3], s_1[g_4], s_1[g_5], s_1[g_6]) \quad (1)$$

	s1	s2	s3
g1	●	●	
g2	(s4: g1)		(s4: -g3, -g5)
g3	(s4: g1) & (s5: -g4, -g6)	(s4: -g2, -g5) & (s5: -g4, -g5)	
g4			(s5: -g3, -g5)
g5	(s4: g1) & (s5: -g4, -g6)		
g6		(s5: -g4, -g5)	(s5: -g3, -g5)

Figure 1: Example BST for the Cancer Class

to be $(1 \wedge 1) \vee (1 \wedge 0) = 1$. Note that \hat{B} will only evaluate the Table 1 **Cancer** samples to *True*.

For a given class set C_i and boolean expression B we can create a **Boolean association rule (BAR)** of the form $B \Rightarrow C_i$. The interpretation of any such BAR, $B \Rightarrow C_i$, is “if $B(s[g_1], \dots, s[g_n])$ evaluates to true for a given sample s , then s should belong to class C_i .” From this point on we will work with the following generalized definitions of support and confidence:

Support: The support of any BAR $B \Rightarrow C_i$, represented as $\text{supp}(B \Rightarrow C_i)$ is:

$$\{\text{samples } s \in C_i \text{ s.t. } B(s[g_1], \dots, s[g_n]) \text{ evaluates to true}\}.$$

The corresponding numerical support value of $B \Rightarrow C_i$ is denoted as $|\text{supp}(B \Rightarrow C_i)|$.

Confidence: The confidence of a BAR $B \Rightarrow C_i$ is

$$\frac{|\text{supp}(B \Rightarrow C_i)|}{|\{\text{samples } s \text{ s.t. } B(s[g_1], \dots, s[g_n]) \text{ evaluates to true}\}|}$$

For CARs these definitions coincide with the CAR definitions of support and confidence found in [2, 3]. Hence, they are natural generalizations of the previous definitions (see section 4.3).

Consider our example boolean expression \hat{B} in terms of Table 1. We can see that the BAR $\hat{B} \Rightarrow \text{Cancer}$ (shown in Eq. 1) has support 3 and confidence 1.

3. BSTS AND BARS

The discussion in this section will focus on tables for each class C_i . These tables, called Boolean Structure Tables (BSTs), will form the basis for our classification method. In order to motivate the utility of BSTs for classification, we will present their close relationship to a special category of BARS which, in turn, will be related back to CARs. Through this discussion we will demonstrate that BSTs contain all the information of the high confidence CARs already known to be valuable for microarray data classification.

3.1 Boolean Structure Tables

A **Boolean Structure Table (BST)** $T(i)$ is a two dimensional table, $T(i) = G \times C_i$, where each table entry refers to a maximum of $|S| - |C_i|$ lists of up to $|G|$ genes each. For every C_i the associated BST, $T(i)$, will require $O((|S| - |C_i|) \cdot |G| \cdot |C_i|)$ space and can be constructed with the same time complexity via Algorithm 1.

When the Algorithm 1 is run on the Table 1 example input and for class **Cancer**, the Boolean Structure Table shown in Figure 1 is produced. In Figure 1 a black dot at location (g, s) indicates that no healthy samples express gene g but some cancerous sample does. A cell (g, s) is left blank only if sample s didn’t express gene g . If (g, s) contains a list of the form $(h : -g_1, \dots, -g_n)$ it means that

Algorithm 1 Create-BST: The BST Creation Algorithm

```

1: Input: Finite set of Genes  $G$ , set of samples  $S$ , Class  $C_i$ 
2: Output: The BST Table for Class  $C_i$ .
3: for all  $(c, h) \in C_i \times S - C_i$  do
4:   initialize a pointer  $\leftarrow$  NULL
5: end for
6: for all  $(g, c) \in G \times C_i$  s.t.  $g \in c$  and  $g \notin \cup_{h \in S - C_i} h$  do
7:   Set  $BST(g, c) \leftarrow$  Black Dot
8: end for
9: for all  $(g, c, h) \in G \times C_i \times S - C_i$  s.t.  $g \in c$  and  $g \in h$  do
10:  if pointer  $(c, h) \neq$  NULL then
11:    push a copy of  $(c, h) \rightarrow BST(g, c)$ 
12:  else
13:     $L = \{g \in G \text{ s.t. } g \in h \ \& \ g \notin c\}$ 
14:    if  $L \neq \emptyset$  then
15:       $(c, h) \leftarrow$   $L$ ’s address
16:    else
17:       $L = \{g \in G \text{ s.t. } g \notin h \ \& \ g \in c\}$ 
18:       $(c, h) \leftarrow$   $L$ ’s address
19:    end if
20:  end if
21:  Push a copy of  $(c, h) \rightarrow BST(g, c)$ .
22: end for

```

s may be distinguished from sample h by the non-expression of any one of genes g_1 through g_n . Similarly, if (g, s) contains a list of the form $(h : g_1, \dots, g_n)$ it means that s may be distinguished from sample h by the expression of any one of genes g_1 through g_n . Such lists will hereafter be referred to as **exclusion lists**.

Note that there is no reason why the BST in Figure 1 was created for the **Cancer** class. We can just as easily build a BST for the **Healthy** class using the example shown in Table 1. In general, if a relational gene expression dataset contains N classes, we can construct N different BSTs for the data set (one for each class).

3.1.1 Runtime Complexity for BST Creation

We can see that the total time to construct BSTs via Algorithm 1 for all of C_1, \dots, C_N is $O\left(\sum_{i=1}^N (|S| - |C_i|) \cdot |C_i| \cdot |G|\right)$. Given that the class sets C_i are all disjoint, we have $\sum_{i=1}^N (|S| - |C_i|) \cdot |C_i| \cdot |G| \leq \sum_{i=1}^N |S| \cdot |C_i| \cdot |G| = |S|^2 \cdot |G|$. Hence, BSTs can be constructed for all C_i s in time $O(|S|^2 \cdot |G|)$.

3.2 BST Generable BARS

We view every BST cell, (g, c) , as an atomic 100% confident BAR. For example, Figure 1’s $(g3, s1)$ -cell corresponds to the BAR

$g3$ expressed *AND* $g1$ expressed *AND* (either $g4$ or $g6$ not expressed) \Rightarrow **Cancer**.

We refer to this rule as the Figure 1 BST’s $(g3, s1)$ -**cell rule**. Note that the cell rule is both (i) 100% confident, and (ii) supported by sample $s1$. Throughout the remainder of this section we will use such cell rules as atomic building blocks to construct more complicated BARS. Furthermore, in Section 5, we will directly employ BST cell rules to build a new classifier called BSTC.

3.2.1 Mining More Complicated BST BARS

Let $T(i)$ be a BST for sample type C_i . We can view each row of $T(i)$ as a 100% confident BAR by combining the row’s cell rules. To see this, choose any $g_j \in G$ and consider the CAR $g_j \Rightarrow C_i$. This CAR can be augmented with exclusion list clauses from each of $T(i)$ ’s g_j -row cells via Algorithm 2. The result will be a BAR with 100% confidence which is logically equivalent to a disjunction of $T(i)$ ’s g_j -row cell rules. See Figure 2 for the gene row BARS

Algorithm 2 BSTRowBAR: Constructing BST Gene Row BAR

```
1: Input: Class  $C_i$ , BST for the class  $T(i)$ , gene  $g_j$ 
2: Output: Row BAR for gene  $g_j$  with 100% conf.
3:  $A \leftarrow FALSE$ 
4: for all  $s \in C_i$  s.t.  $T(i)$ 's  $(g_j, s)$  - cell is not empty do
5:    $B \leftarrow TRUE$ 
6:   for all exclusion lists  $e \in T(i)$ 's  $(g_j, s)$ -cell do
7:     if  $e = (s_k : -g_{l_1} \dots -g_{l_m})$  then
8:        $B \leftarrow B \text{ AND } (-g_{l_1} \text{ OR } \dots \text{ OR } -g_{l_m})$ 
9:     else if  $e = (s_k : g_{l_1} \dots g_{l_m})$  then
10:       $B \leftarrow B \text{ AND } (g_{l_1} \text{ OR } \dots \text{ OR } g_{l_m})$ 
11:     end if
12:    $A \leftarrow A \text{ OR } B$ 
13: end for
14: end for
15: Return  $g_j \text{ AND } A \Rightarrow C_i$ 
```

Gene g_1 : (g_1 expressed) \Rightarrow Cancer.

Gene g_2 : (g_2 expressed AND [EITHER (g_1 expressed) OR (either g_5 or g_3 not expressed)]) \Rightarrow Cancer.

Gene g_3 : (g_3 expressed AND [EITHER $\{(g_1$ expressed) AND (either g_4 or g_6 not expressed)} OR $\{($ either g_2 or g_5 not expressed) AND (either g_4 or g_5 not expressed)}]]) \Rightarrow Cancer.

Gene g_4 : (g_4 expressed AND [either g_5 or g_3 not expressed]) \Rightarrow Cancer.

Gene g_5 : (g_5 expressed AND [g_1 expressed AND (either g_4 or g_6 not expressed)]) \Rightarrow Cancer.

Gene g_6 : (g_6 expressed AND [(either g_4 or g_5 not expressed) OR (either g_3 or g_5 not expressed)]) \Rightarrow Cancer.

Figure 2: Gene Row BARs with 100% Confidence Values.

which result from applying Algorithm 2 to the example BST in Figure 1.

For the remainder of this paper we will restrict our attention to BARs that may be generated by taking conjunctions of BST cell rule disjuncts. Henceforth we simply refer to these as BARs. It is very important to notice that all such BARs have a special form: Their antecedents consist of a CAR antecedent ANDed with a disjunction of BST exclusion list clause conjunctions. Consider the BAR for gene g_6 in Figure 2. Gene g_6 's rule antecedent consists of a CAR antecedent, g_6 , conjoined to a disjunction of the Figure 1 exclusion list clauses: (either g_4 or g_5 not expressed) and (either g_3 or g_5 not expressed).

Along these same lines, BARs with more complex antecedents can be created by taking the logical AND of a BST's gene row rules. For example, consider our running example BST's gene row rules listed in Figure 2. We can form the 100% confident CAR (g_1 expressed AND g_6 expressed) \Rightarrow Cancer by ANDing Figure 2 gene row rules for g_1 and g_6 as follows: While ANDing we use the BST in Figure 1 to quickly simplify the resulting expression. First, we can tell that product will only be supported by sample s_2 because only the BST's s_2 column contains non-empty cells for both of gene rows g_1 and g_6 . Thus, we only need to consider the exclusion lists in cells (g_1, s_2) and (g_6, s_2) while forming our product. Second, the black dot in BST entry (g_1, s_2) means we don't have to use the Healthy sample s_5 exclusion list information $(s_5 : -g_4, -g_5)$ from BST entry (g_6, s_2) in our new rule. This is because gene g_1 already excludes s_5 on its own since $g_1 \notin s_5$. By ANDing gene row rules in this manner we can create BARs with antecedents that are the conjunction of any desired CAR antecedent with a simplified exclusion list based clause (to eliminate non- C_i supporting samples).

Algorithm 3 (MC)²BAR Mining Algorithm

```
1: Input: Number MC2BARs  $k$ , Class  $C_i$ , BST for the class  $T(i)$ , Genes  $G$ 
2: Output: Set of Top- $k$  supported MC2BARs
3: for all  $g \in G$  do
4:   initialize  $C_i\_SUP[g] \leftarrow C_i$  subset supporting  $T(i)$ 's 100% confidence  $g$ -row BAR
5: end for
6: Sort  $C_i\_SUP$  by subset size
7: Create arrays  $RULES$  and  $RULE\_SUP$ 
8:  $B \leftarrow$  largest size of any  $C_i\_SUP$  subset.
9: for all  $S \in C_i\_SUP$  s.t.  $|S| = B$  do
10:   $R \leftarrow$  AND all  $T(i)$  gene-row rules with support  $\supseteq S$ 
11:  Push  $R \rightarrow RULES$ 
12:  Push  $S \rightarrow RULE\_SUP$ 
13:  Push  $S \rightarrow B\_SUP$ 
14: end for
15: for all  $(S_1, S_2) \in B\_SUP \times RULE\_SUP$  do
16:  Push  $S_1 \cap S_2 \rightarrow NEWSUPP$ 
17:  Sort  $NEWSUPP$  by subset size
18:  Merge  $NEWSUPP$  into  $C_i\_SUP$  (ordered)
19:  Remove duplicates from  $C_i\_SUP$ 
20: end for
21: Set  $C_i\_SUP = C_i\_SUP - B\_SUP$ 
22: Empty  $B\_SUP$ 
23: if  $|RULES| < k$  then
24:   Goto step 8
25: else
26:   Quit.
27: end if
```

4. EXTENDED CAR MINING SECTION

We begin by introducing the concept of maximally complex BARs which we need in order to both accomplish fast BAR mining via BSTs and to better understand the relationship between CARs and BARs.

4.1 Maximally Complex BARs

We'll say a BAR is **maximally complex** if no new genes may be conjoined to its antecedent's CAR portion without decreasing the size of its C_i support set. For example, by looking at Figure 1 above we can see that the Figure 2 BARs for genes g_2 and g_6 are maximally complex while all others are not. The g_2 -row rule has support sample set $\{s_1, s_3\}$ which isn't a subset of any other gene row rule's support set. Likewise, the g_6 -row rule has support that isn't a subset of any other gene row rule's support. However, the g_1 -row rule (for example) isn't maximally complex because its support $\{s_1, s_2\}$ is the same as the support of the g_3 -row rule.

With the concept of maximum complexity in hand we can extend the idea of gene-row BAR generation to find a **Maximally Complex 100%** (i.e. **Maximally**) **Confident Boolean Association Rule**, or **(MC)²BAR**, for each of the largest k supportable C_i sample subsets. Toward this end we propose the Mine-MCMC BAR algorithm (shown in Algorithm 3). This algorithm finds a (MC)²BAR for each of the top- k supporting C_i sample subsets. Note that Mine-MCMC BAR algorithm can be executed in worst case time $O(k^2 \log(k) \cdot |G| \log(|G|) \cdot |S|^2)$. Furthermore, the algorithm is *progressive*. New BARs can be shared with the user as they are added to the $RULES$ array.

Upon completion, the algorithm will store a (MC)²BAR for each top- k supporting C_i subset in the $RULES$ array. It is important to observe that the more complex a boolean association rule is (i.e. the more ANDed gene row BARs used to create it) the fewer exclusion list clauses it must contain. In general, if R_1 and R_2 are any two Boolean association rule antecedents, then $(R_1 \text{ AND}$

Algorithm 4 Mine-MCMCBBAR-Per-Samp Mining Algorithm

1: **Input:** Number MC²BARs k , Class C_i , BST for the class $T(i)$, Genes G
2: **Output:** Set of Top- k supported MC²BARs per Sample
3: Initialize $ALL_RULES \leftarrow \emptyset$
4: **for all** $c \in C_i$ **do**
5: Find k rules with support containing c via modified Mine-MCMCBBAR
6: Push k rules $\rightarrow RULES$
7: Sort $RULES$ by support (lexicographically)
8: Merge $RULES$ with ALL_RULES
9: Remove duplicates from ALL_RULES
10: Empty $RULES$
11: **end for**

$R_2) \Rightarrow C_i$ must only contain exclusion clauses for samples excluded in both R_1 and R_2 .

Note that in the Mine-MCMCBBAR algorithm a secondary ordering relation could be used to break ties between same-sized subsets in the $C_i_SUPPORTS$ array. For example, subsets whose related (MC)²BARs have a small number of excluded samples could be ordered before subsets of the same size whose related (MC)²BARs have many excluded samples. This would effectively break ties by ordering (MC)²BARs with higher confidence antecedent CAR portions first (see section 4.3).

We next turn our attention to the correctness of the *Mine – MCMCBBAR* algorithm. We can formally prove that the *Mine – MCMCBBAR* algorithm is guaranteed to produce a top- k set of MC²BARs, of maximal support.

THEOREM 1. *The Mine-MCMCBBAR algorithm will produce a top- k set of (MC)²BARs of maximal support.*

PROOF:

Let g be a gene. $T(i)$'s g -row BAR will have support \supseteq the support of any (MC)²BAR whose antecedent's CAR portion contains g . More generally, ANDing two (MC)²BARs will always yield a BAR whose support is the intersection of its parents' supports. Below we will refer to the number of CAR antecedent genes in a BAR as that rules *complexity*. A *simple* BAR has only one CAR antecedent gene (e.g., BST gene-row BARs and atomic cell rules are simple). Let $B \Rightarrow C_i$ be any (MC)²BAR implying class C_i . Then, there exists a set of (MC)²BARs $B_1 \Rightarrow C_i, \dots, B_n \Rightarrow C_i$ such that the following are true:

1. The antecedents B_1, \dots, B_n are all conjunctions of antecedents of 100% confident maximally supported gene row rules (as discussed in section 3.2 above).
2. The subset of C_i samples which evaluate the antecedent B to true is exactly the subset of C_i samples which evaluate the conjunction of the antecedents B_1, \dots, B_n to true.
3. For each $B_j, 1 \leq j \leq n$, the subset of C_i samples which evaluate the antecedent B_j to true is a superset of the subset of C_i samples which evaluate the antecedent B to true.
4. If $B \Rightarrow C_i$ is non-simple, then for each non-simple rule antecedent $B_j, 1 \leq j \leq n$, there exist two 100% confident simple rule antecedent clauses, x and y , such that $B_j = (x \text{ AND } y) \text{ AND } R$. Here $R \Rightarrow C_i$ is either a (MC)²BAR or true. Furthermore, the subset of C_i samples which evaluate the antecedent B_j to true is exactly the subset of C_i samples which evaluate $(x \text{ AND } y)$ to true.

By induction on the number of non-exclusion clause genes in B (i.e. the degree of complexity) we can see that the 4 part claim above is true. Using the claim we can see that if we visit each possible supporting C_i set from largest to smallest, we will be able use the BST to generate a (MC)²BAR for each visited supporting set. Using this claim we can see that the *Mine – MCMCBBAR* algorithm above will correctly mine the top- k supported (MC)²BARs (one for each possible supporting set). \square

It is often desirable to make sure that each training sample belongs to the support of at least one mined rule. The Mine-MCMCBBAR algorithm can easily be modified to accomplish this task (i.e. to find a (MC)²BAR for each top- k supportable C_i subset containing any given C_i sample c). All we have to do is restrict our attention to C_i sample subsets containing sample c in the Mine-MCMCBBAR algorithm. Hence, if desired we may ensure that for every C_i sample c we generate a (MC)²BAR for each of the top- k supportable C_i subsets containing c . The Mine-MCMCBBAR-Per-Samp algorithm shown in Algorithm 4 does this in worst case time complexity $O(k^2 \log(k) \cdot |G| \log(|G|) \cdot |S|^3)$.

Note that (MC)²BARs are analogous to the rule group upper bounds proposed in [10]. Hence, they can be used in the same way to represent all BST creatable BARs with the same support set.

4.2 Interesting Boolean Rule Groups

The following rule group definitions are essentially identical to those found in FARMER and similar results (such as the uniqueness of rule group upper bounds [10]) are true. An **interesting boolean rule group (IBRG)** is as follows:

IBRG Definition 1. Let S be the set of samples in our data set. We say that $RG = \{R_I \Rightarrow C_i | R_I \text{ is a conjunction of simple 100% confident BAR antecedents}\}$ is an **IBRG** with consequent C_i and antecedent support set S if it has both the following properties:

1. $\forall R_I \Rightarrow C_i \in RG$ we have $supp(R_I \Rightarrow C_i) = S$.
2. All conjunctions of simple 100% confident BAR antecedents, R_I , with row support set S have $R_I \Rightarrow C_i \in RG$.

A 100% confident BAR $R_u \Rightarrow C_i$ in RG is called an **upper bound** of $RG \iff$ there exist no $\hat{R} \Rightarrow C_i \in RG$ such that \hat{R} 's CAR portion $\supset R_u$'s CAR portion. A 100% confident BAR $R_L \Rightarrow C_i \in RG$ is called a **lower bound** of $RG \iff$ there exist no $\hat{R} \Rightarrow C_i \in RG$ such that \hat{R} 's CAR portion $\subset R_L$'s CAR portion.

In our running example the set $\{(g_1 \text{ AND } g_6 \text{ expressed}) \Rightarrow Cancer, (g_3 \text{ AND } g_6 \text{ expressed AND (either } g_4 \text{ or } g_5 \text{ not expressed)}) \Rightarrow Cancer, (g_1 \text{ AND } g_3 \text{ AND } g_6 \text{ expressed}) \Rightarrow Cancer\}$ is the boolean rule group with consequent *Cancer* and antecedent support set $\{s2\}$. The first two boolean association rules are the lower bounds while the last most complex rule is the upper bound. We can see that each maximally complex rule found by the two algorithms in section 4.1 will be an upper bound to a unique IBRG.

Note that we are now in the position to use (MC)²BARs to classify query gene expression data samples. To do so we can (i) Mine the top- k supported IBRG upper bounds per training sample (via Figure 4 algorithm) for each class label/type in the data, (ii) Calculate a query classification number $\in [0, 1]$ for every generated IBRG upper bound by using each BAR's exclusion lists (see section 5.2), and then (iii) Classify the query sample as the class type antecedent belonging to the IBRG upper bound with the largest classification number. This approach is polynomial time. However, it also depends on the support related parameter k . Hence, we forgo

greater development of these classification schemes and concentrate on the support pruning/parameter free classification method presented in section 5.3.

Next we formalize the relationship between BST created boolean association rules and conjunctive association rules.

4.3 BARs Relationships to CARs

Let $R \Rightarrow C_i$ be any 100% confident BST created BAR containing exclusion clauses for non- C_i samples h_1, \dots, h_m . Removing all exclusion list clauses related to $\{\hat{h}_1, \dots, \hat{h}_p\} \subset \{h_1, \dots, h_m\}$, $p \leq m$, will create a new boolean association rule, $\hat{R} \Rightarrow C_i$, with support = $\text{supp}(R \Rightarrow C_i)$ and confidence $\geq \frac{|\text{supp}(R \Rightarrow C_i)|}{|\text{supp}(R \Rightarrow C_i)| + p}$. Let's consider the g_3 -row BAR from our running example:

(g_3 expressed AND [EITHER $\{(g_1 \text{ expressed}) \text{ AND } (\text{either } g_4 \text{ or } g_6 \text{ not expressed})\}$ OR $\{(\text{either } g_2 \text{ or } g_5 \text{ not expressed}) \text{ AND } (\text{either } g_4 \text{ or } g_5 \text{ not expressed})\}$]) \Rightarrow Cancer.

It has 100% confidence and support $\{s1, s2\}$. Now, if we remove all exclusion list clauses related to sample row $s5$ we end up with the boolean association rule (g_3 expressed AND [EITHER (g_1 expressed) OR (either g_2 or g_5 not expressed)]) \Rightarrow Cancer. This new rule has support $\{s1, s2\}$ and a confidence of $\frac{|\{s1, s2\}|}{|\{s1, s2, s5\}|} = \frac{2}{3}$. The preceding observation leads us to the following theorem:

THEOREM 2. *Let D be a relational data set containing s samples no two of which are the same (i.e. no two sample rows express the exact same set of genes). Then, there exists a pure conjunction B implying a class type C (i.e., a CAR) with confidence c and support supp for $D \iff$ there exists a 100% confident BST generated BAR $\hat{B} \Rightarrow C$ for D that: (i) has $\text{supp}(\hat{B} \Rightarrow C) = \text{supp}$, and (ii) contains exclusion list clauses actively excluding $(\frac{1}{c} - 1)|\text{supp}|$ non- C samples.*

PROOF. \Leftarrow : From the observation directly preceding this theorem we can see that if $\hat{B} \Rightarrow C$ has $\text{supp}(\hat{B} \Rightarrow C) = \text{supp}$ then removing all the exclusion list clauses from \hat{B} (by replacing them all with true) will create a new pure conjunction B with $\text{supp}(B \Rightarrow C) = \text{supp}$. Furthermore, we require that $\{\text{non-}C \text{ samples excluded by exclusion clauses}\} = \{\text{non-}C \text{ samples satisfying } B\}$ (i.e., the exclusion clauses actually exclude something). Hence, $B \Rightarrow C$ will have confidence $c = \frac{|\text{supp}|}{|\text{supp}| + \#\text{excluded samples}}$. \Rightarrow : Let B be a conjunction of items/genes g_1, \dots, g_n . Given that no two samples in D are the same we can build a 100% confident BST for class C of D . Furthermore, both the following are true:

1. A non- C sample h expresses all genes $g_1, \dots, g_n \iff \forall s \in \text{supp}$ and $1 \leq i \leq n$ the BST cell (g_i, s) contains an active exclusion list for h . Thus, only non- C samples expressing all of g_1, \dots, g_n (and therefore satisfying B) generate active exclusion lists in all relevant (g_i, s) BST cells.
2. $\text{supp}(B \Rightarrow C) = \bigcap_{1 \leq j \leq n} \text{supp}(g_j \Rightarrow C)$.

Here we get the \hat{B} by ANDing down each of the BST $\text{supp}(B \Rightarrow C)$ sample column's g_i cells and then ORing the resulting $|\text{supp}(B \Rightarrow C)|$ rules together. \square

Theorem 2 tells us how we can get CARs from BARs. Furthermore, it says 100% confident BARs with large support and a small number of excluded samples are equivalent to high support/confidence CARs. Hence, genes that show up in many high confidence, high support CARs will also be prevalent in many 100% confident BARS with high support and a low number excluded

samples. Most importantly, we see that all high confidence CARs (which tend to be good classifiers) have closely related BAR counterparts. Furthermore, these counterparts can be mined from a BST by ANDing gene row BARs.

5. BST-BASED CLASSIFICATION

In principal, 100% confident BST-generable BARs should be sufficient for classification because they contain at least as much information as all generable CARs do (see section 4.3). Indeed, beyond what CARs with similar support are capable of, 100% confident BARs supply us with "unpolluted" ground truth. Thus, it's not too surprising that the class of BST-generable BARs we've looked at so far will be enough to enable highly accurate classification.

Let C_i be a class set of interest and $T(i)$ be the BST for class C_i constructed from the given training data. From section 3.2 we can see that all BST generable BARs for class C_i are created by combining $T(i)$ cell rules. Thus, we expect that by restricting our attention to the $O(|G| \cdot C(i))$ atomic $T(i)$ cell rules we will be, in some sense, still considering all $T(i)$ generable BARs for C_i . Our new scalable classifier, the **Boolean Structure Table Classifier (BSTC)**, capitalizes on this line of thought by ignoring BAR generation and focusing exclusively on atomic BST cell rules.

5.1 BSTC Overview

Let Q be a test/query gene expression data sample and $T(i)$ be a BST for class set C_i . BSTC is a heuristic rule-based classifier motivated by standard Boolean formula arithmetization techniques [23] such as those employed in fuzzy satisfiability [29]. By using these ideas we can avoid the highly costly process of support/confidence based association rule mining. Instead of explicitly generating rules, BSTC decides (heuristically), for all C_i , how well Q collectively satisfies $T(i)$'s atomic cell rules. BSTC then classifies Q as the sample class whose BST has the highest expected atomic rule satisfaction level from Q .

Intuitively, we expect BSTC to be accurate because it approximates the results of CAR-based classification: Suppose that a high support/confidence CAR exists which classifies our query sample Q as class C_j . This will only happen if all the CAR's antecedent genes, AG , appear in both (i) Q and, (ii) most of the training samples in the CAR's consequent class C_j . Let $T(j)$ be the BST for class C_j . Because of (ii) most of $T(j)$'s sample columns must contain cell entries for all the AG genes. Furthermore, all $T(j)$'s AG cell entries will have few exclusion lists in common (by Theorem 1). Hence, $T(j)$'s expected atomic rule satisfaction level from Q (i.e., Q 's classification value) should be heavily influenced (increased) by the AG rows and their few shared lists.

5.2 BST Cell Rule Satisfaction

As above, let Q be a test/query gene expression data sample and $T(i)$ be a BST for class set C_i . Algorithm 5, BSTCE, gives BSTC's method of calculating the level that Q satisfies a given atomic $T(i)$ cell rule. We next explain the rational behind BSTCE.

We know that each $T(i)$ (g, s)-cell exclusion list, L , corresponds to a disjunction in $T(i)$'s (g, s)-cell rule. Hence, if Q satisfies any one negation/inclusion in L , Q will satisfy L . However, if Q expresses most of its genes in common with L 's associated non- C_i sample we assume it's probably not of type C_i (i.e., Q is weakly excluded). Hence, we use BSTCE's line 4 ratio to approximate the probability that L correctly excludes Q from being of L 's associated sample's class.

In order for the (g, s)-cell rule to be satisfied, all of (g, s)'s exclusion lists must be satisfied (i.e., logical AND). If independence of each exclusion list's correct classification is assumed it's natu-

Algorithm 5 BST Cell rule quantized Evaluation (BSTCE)

1: **Input:** Class C_i , BST for the class $T(i)$, Samples S , Query sample Q
2: **Output:** Classification value
3: **for all** non-empty exclusion lists e in $T(i)$'s cells **do**
4: $V_e \leftarrow \frac{|\{\hat{g} \in e \text{ s.t. } Q[\hat{g}] = 1\}|}{|e|}$
5: **end for**
6: **for all** $(g, s) \in \{g \in G \text{ s.t. } Q[g] = 1\} \times C_i$ **do**
7: **if** $T(i)(g, s)$ contains a \bullet **then**
8: $T(i)[g][s] \leftarrow 1$
9: **else**
10: $T(i)[g][s] \leftarrow \text{Min} \{V_e \text{ s.t. } e \text{ is in } T(i)(g, s)\}$
11: **end if**
12: **end for**
13: **for all** non-blank sample columns $s \in T(i)$ **do**
14: $V_s \leftarrow \text{Mean of non-blank } T(i)[*][s] \text{ values}$
15: **end for**
16: Return the Mean of step 16's V_s values

Algorithm 6 The BSTC Algorithm

1: **Input:** BSTs for all dataset classes $T(1), \dots, T(N)$, Query sample Q
2: **Output:** Classification for query sample Q
3: **for all** $i \in \{1, \dots, N\}$ **do**
4: $CV(i) \leftarrow \text{BSTCE}(T(i), Q)$
5: **end for**
6: Return $\min\{i | CV(i) = \max\{CV(1), \dots, CV(N)\}\}$

ral to multiply all of (g, s) 's list's probabilities. We don't assume independence and use a min instead (line 10).

Finally, recall that all black dots in $T(i)$ correspond to genes expressed only in class C_i samples. If Q expresses a black dot gene it automatically satisfies all that gene's non-empty $T(i)$ cell rules. Hence, black dots are all assigned values of 1 in BSTCE's line 8.

Once we have used BSTCE lines 1-12 to calculate Q 's classification values (i.e., $T(i)$'s atomic rule satisfaction levels from Q) for each relevant simple (g, s) -cell rule, we are nearly finished. We have all the values required to judge Q 's similarity to $T(i)$ via an expectation calculation. For the sake of $T(i)$'s expectation calculation, all that is left to do is imagine choosing a relevant simple $T(i)$ rule at random and then using it to classify Q . To randomly select a (g, s) rule we first imagine selecting a non-empty $T(i)$ sample column uniformly at random and then picking a cell-rule from that column uniformly at random. The expected probability of correctly classifying Q with $T(i)$ via this method (which heuristically is proportional to $T(i)$'s expected satisfaction level from Q) is then calculated by averaging the approximate cell rule satisfaction levels down each non-empty sample column (line 14) and then averaging the resulting non-empty sample averages (line 16).

5.3 BSTC Algorithm

Suppose we are given relational training data D containing sample rows S split up into disjoint class sets C_1, \dots, C_N . BSTC uses D to construct N BSTs, $T(1), \dots, T(N)$. Next, let G be the union of the elements contained in each sample row of D (i.e. the gene set of D) and let Q be a query sample with expression information regarding G . BSTC will use the BSTCE algorithm to classify Q as being the C_i with smallest i such that $\text{BSTCE}(T(i), Q) = \max\{\text{BSTCE}(T(j), Q) | 0 \leq j \leq N\}$. See Algorithm 6 for the BSTC algorithm.

Note that there is no reason why N must be 2. *BSTC easily generalizes to datasets containing more than two class labels.*

5.3.1 BSTC Runtime

As noted in section 3.1.1 it takes time and space $O(|S|^2 \cdot |G|)$

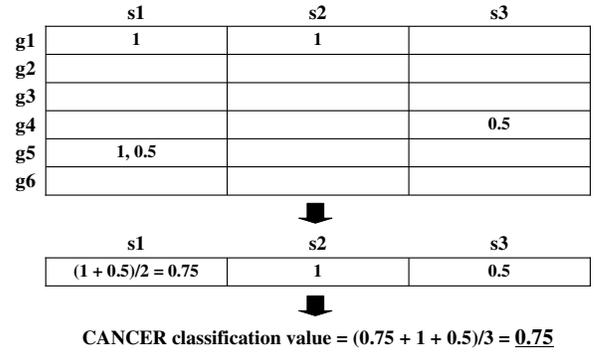


Figure 3: BSTC cell rule Evaluation Example

to construct all the BSTs $T(1), \dots, T(N)$. Thus, BSTC requires time and space $O(|S|^2 \cdot |G|)$ to construct. Furthermore, during classification BSTC must calculate $\text{BSTCE}(T(i), Q)$ for $1 \leq i \leq N$. BSTCE (Algorithm 5) runs in $O((|S| - |C_i|) \cdot |G| \cdot |C_i|)$ time per query sample. Therefore the BSTC worst case evaluation time is also $O(|S|^2 \cdot |G|)$ per query sample. See Section 8 for more on BSTC's per-query classification time.

5.3.2 Biological Meaning of BSTC Classification

Association rules mined from gene expression data provide an intuitive representation of biological knowledge (i.e., the expression of certain genes implies cancer). Hence, CAR-based classifiers have the desirable ability to justify each non-default consequent class query classification with the biologically meaningful CAR(s) the query satisfied. BSTC, being rule-based and explicitly related to CAR-classifiers, also has this property.

BSTC can support its query classifications with BARs of any user specified complexity. Most simply, for any given query sample Q and $c \in (0, 1]$, BSTC can justify its classification of Q as class C_i by reporting all $T(i)$ atomic cell rules with satisfaction levels $\geq c$. Note that returning this information requires no additional per-query classification time. Also note that section 3.2.1 methods can be used to progressively mine more complex highly satisfied BARs if desired.

5.4 BSTC Example

Consider our running example from Table 1. In order to construct BSTC we must construct both $T(\text{Healthy})$ and $T(\text{Cancer})$ (shown in Figure 1). Once both BSTs have been constructed we can begin to classify query samples. Suppose, for example, we are given the query sample $Q = \{g_1 \text{ expressed, } g_2 \text{ not expressed, } g_3 \text{ not expressed, } g_4 \text{ expressed, } g_5 \text{ expressed, } g_6 \text{ not expressed}\}$. To classify this query we must first calculate $\text{BSTCE}(T(\text{Cancer}), Q)$ and $\text{BSTCE}(T(\text{Healthy}), Q)$.

The evaluation of $\text{BSTCE}(T(\text{Cancer}), Q)$ proceeds as follows: Since our query sample Q expresses gene g_5 we can see that we must, for example, determine the fraction of both of the (g_5, s_1) -cell's exclusion lists satisfied by Q . The (g_5, s_1) -cell's $(s_4 : g_1)$ exclusion list is totally satisfied since Q expresses g_1 . Hence, it gets a value of 1. However, the $(s_5 : -g_4, -g_6)$ exclusion list is only half satisfied since, although Q doesn't express g_6 , Q does express g_4 . Thus, in total, we only consider half of the simple (g_5, s_1) -cell rule to be satisfied (i.e. the s_5 exclusion list is the weakest link). Continuing to use BSTC's approximation scheme

for the expected probability of Q 's correct Cancer classification via the Figure 1 BST we obtain Figure 3. Note that only Figure 3 gene rows corresponding to genes expressed in Q are non-empty.

If we now evaluate $\text{BST-EXPECT}(T(\text{Healthy}), Q)$ we obtain a final value of $\frac{3}{8}$. To finish, BSTC will compare Q 's Cancer classification value of $\frac{3}{4}$ to Q 's Healthy classification value of $\frac{3}{8}$ and conclude that Q is most probably Cancer. Hence, Q will be classified as Cancer.

6. EXPERIMENTAL EVALUATION

All experiments reported here were carried out on a 3.6 GHz Xeon machine with 3GB of memory running Red Hat Linux Enterprise 4. For our empirical evaluation we use four standard real microarray datasets¹. Table 2 lists the dataset names, class labels, and the number of samples of each class. All discretization was done using the entropy-minimized partition².

Dataset	# Genes	Class 1 label	Class 0 label	# Class 1 samples	# Class 0 samples
ALL/AML (ALL)	7129	ALL	AML	47	25
Lung Cancer (LC)	12533	MPM	ADCA	31	150
Prostate Cancer (PC)	12600	tumor	normal	77	59
Ovarian Cancer (OC)	15154	tumor	normal	162	91

Table 2: Gene Expression Datasets

Executables for both RCBT and Top-k were provided by the authors of [9]. In all experiments, the Top-k rule generator was used to generate rule groups for RCBT. Unless otherwise noted we ran both Top-k and RCBT with the author suggested parameter values (i.e., support = 0.7, $k = 10$, $nl = 20$, 10 RCBT classifiers). Hence, while generating rules for RCBT we used Top-k with a minimum support value of 0.7 and found the 10 most confident covering rule groups (i.e. $k = 10$). Furthermore, during classification we used RCBT with the suggested 10 classifiers (1 primary and 9 standby). Finally, nl , the number of lower bound rules to use for classification per Top-k mined rule group, was set equal to 20. We coded BSTC with C++³.

6.1 Preliminary Experiments

Each of Table 2's four gene expression datasets comes with a clinically determined training set. The authors of [9] provided us with their discretizations of these four datasets. We ran BSTC on their discretizations and BSTC matched RCBT's reported mean accuracy (about 96%) outperforming CBA (87%), IRG (81%), Weka 3.2 (C4.5 family single tree (74%), bagging (78%), boosting(74%)), and SVM^{light} 5.0 (93%) in reported mean performance [9].

To compare BSTC and RCBT with the most recent R e1071 package SVM implementation [7] and randomForest version 4.5 [6] we discretized the four datasets and reran BSTC/RCBT. To keep comparisons fair we ran SVM and randomForest on the same genes selected by our entropy discretization except with their original undiscretized gene expression values. SVM was run with its default radial kernel. We ran randomForest 10 times with its default 500 trees for ALL, LC, and OC and its accuracy was constant. For

¹All four undiscretized gene expression data files are available at <http://sdmc.i2r.a-star.edu.sg/rp/>

²The entropy partition code is available at http://cran.r-project.org/src/contrib/dprep_1.0.tar.gz

³BSTC Code and data files are located at <http://www-personal.umich.edu/~markiwen/>.

PC we had to increase randomForest's number of trees to 1000 before its accuracy stabilized over the 10 runs.

Table 3 contains the number of class 0/1 samples in the clinically determined training set, the number of genes selected by our entropy discretization, and our experimental results. As shown in this table, the overall average accuracies of BSTC and RCBT are again best at about 96% each. When compared against RCBT, SVM, and randomForest on the individual tests we can see that BSTC is alone in having 100% accuracy on the majority of datasets.

However, BSTC's performance on the preliminary AML/ALL dataset test is relatively poor. This is likely due to over fitting. Every error BSTC made mistook a class 0 (AML) test sample for a class 1 (ALL) test sample (i.e., all errors were made in this same direction). And, the ALL training data has both (i) about 2.5 times as many class 1 samples as class 0 samples, and (ii) a small number of total samples/genes. When the training set is more balanced and the number of samples/genes is larger we can expect that cancellation of errors will tend to neutralize/balance any over fitting effects in BSTC. And, BSTC is a method meant primarily for large training sets where CAR-mining is prohibitively expensive. As we will see later in Section 6.2.1, BSTC's performance is much better for larger AML/ALL training set sizes.

6.2 Cross-Validation Studies

Cross-validation studies make comparisons less susceptible to the choice of a single training dataset and provide performance evaluations that are likely to better represent program behavior in practice. We next present results from a thorough cross-validation study completed using 100 different training/test sets from each of the ALL, LC, PC, and OC data sets. For these cross-validation tests we benchmark BSTC against Top-k/RCBT because (i) BSTC/RCBT perform best in our preliminary experiments, (ii) Top-k/RCBT is the fastest/most accurate CAR-based classifier for microarray data, and (iii) we are interested in BSTC's CAR-related vs Top-k/RCBT's CAR-based scalability.

For the cross-validation study we generated training sets of sizes 40%, 60%, and 80% of the total samples. Each training set was produced by randomly selecting samples from the original combined dataset. We then used the standard R dprep package's entropy minimized partition [1] to discretize the selected training samples. Finally, the remaining dataset samples were used for testing the two classifiers after rule/BST generation on the randomly selected training data. For each training set size we produced 25 independent tests. In addition to these training sets, we created an additional 25 1- x /0- y tests. To create these tests we chose training data by randomly selecting x class 1 samples and y class 0 samples to be used as training data. As before, the remaining samples were then used to test both classifiers. For each dataset the x and y values are chosen so that the resulting 25 classification tests have the exact same training/test data proportions as the single related dataset test reported in section 6.1. For each training set size we plot our results using a boxplot.

Boxplot Interpretation: Each boxplot that we show in this section can be interpreted as follows: The median of the measurements is shown as a diamond, and a box with boundaries is drawn at the first and the third quartile. The range between these two quartiles is called the inter-quartile range (IQR). Vertical lines (a.k.a. "whiskers") are drawn from the bottom and the edge of the box to indicate the minimum and the maximum value, unless outliers are present. If outliers are present, the whiskers only extend to $1.5 \times IQR$. The outliers that are near (i.e. within $3 \times IQR$) are drawn as an empty circle, and outliers that are further are drawn using an asterisk.

Dataset	# Class 1 Training Samples	# Class 0 Training Samples	Genes After Discretization	BSTC Accuracy	RCBT Accuracy	SVM Accuracy	randomForest Accuracy
ALL/AML (ALL)	27	11	866	82.35%	91.18%	91.18%	85.29%
Lung Cancer (LC)	16	16	2173	100%	97.99%	93.29%	99.33%
Prostate Cancer (PC)	52	50	1554	100%	97.06%	73.53%	73.53%
Ovarian Cancer (OC)	133	77	5769	100%	97.67%	100%	100%
Average Accuracy				95.59%	95.98%	89.5%	89.54%

Table 3: Results Using Given Training Data

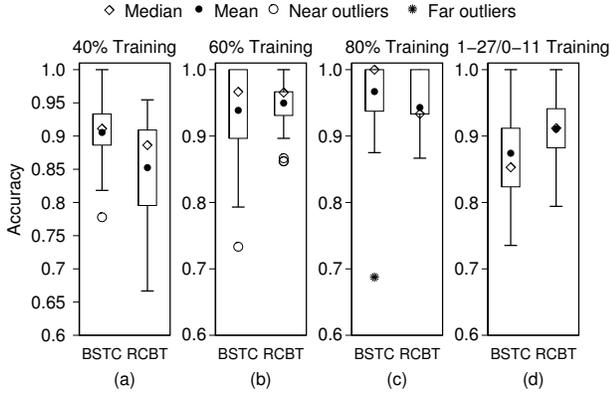


Figure 4: ALL Cross-Validation Results

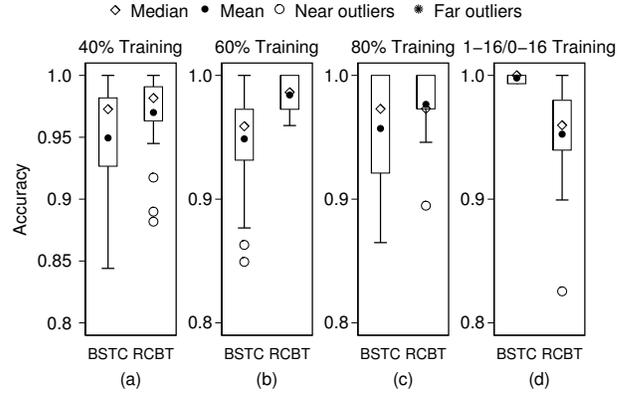


Figure 5: LC Cross-Validation Results

6.2.1 ALL/AML (ALL) Experiment

Figure 4 shows the classification accuracy for the ALL/AML dataset. As can be seen in this figure, BSTC and RCBT have similar accuracy across the ALL/AML tests as a whole. BSTC outperforms RCBT in terms of median and mean accuracy on the 40% and 80% training set sizes while RCBT has better median/mean accuracy on the 1-27/0-11 training size tests. And, both classifiers have the same median on the 60% training set size. Over the 100 ALL/AML tests we see that BSTC has a mean accuracy of 92.13% while RCBT has a mean accuracy of 91.39% (they are very close).

It's noteworthy that BSTC is 100% accurate on the majority of 80% training size tests. However, BSTC appears to have slightly higher variance than RCBT on all but the 40% training tests. Considering all the results together both BSTC and RCBT have essentially equivalent classification accuracies on the ALL/AML dataset.

6.2.2 Lung Cancer (LC) Experiment

The results for the Lung Cancer dataset are reported in Figure 5. Here, again, both BSTC and RCBT have similar classification behavior. RCBT has higher mean and median accuracies on the 40% and 60% tests while BSTC outperforms RCBT on the 1-16/0-16 tests. Meanwhile, both classifier have the same median on the 80% training test. Over all 100 LC tests we find that BSTC has a mean accuracy of 96.32% while RCBT has a mean accuracy of 97.08% (again, they are very close).

As before, BSTC is alone in having 100% accuracy more than half the time for any training set size (see Figure 5 (d)). However, RCBT has smaller variance for 3 of the 4 training set sizes. Therefore, as for the ALL/AML data set, both BSTC and RCBT have essentially the same classification accuracy on LC.

6.2.3 Prostate Cancer (PC) Experiment

RCBT begins to run into a computational difficulties on PC's larger training set sizes. This is because before using a Top-k rule group for classification RCBT must first mine nl lower bound rules for the rule group. RCBT accomplishes rule group lower bound mining via a pruned breadth-first search on the subset space of the rule group's upper bound antecedent genes. This breadth-first search can be quite time consuming. In the case of the Prostate Cancer (PC) dataset all 100 classification tests (25 tests for each of the 4 training set sizes) generated at least one top-10 rule group upper bound with more than 400 antecedent genes. Due to the difficulties involved with a breadth-first search over the subset space of a several hundred element set, RCBT began suffering from long run times on many PC classification tests.

Table 4 contains four average classification test run times (in seconds) for each PC training size. The 'BSTC' column run times reflect the average time required to build both class 0 and class 1 BSTs and then use them to classify all the test samples. Each 'Top-k' column run time is the average time required for Top-k to mine the top 10 covering rule groups (with minimum support 0.7) for each training set.

Table 4's 'RCBT' column gives average run times for RCBT using a time cutoff value of 2 hours for all the training sets. For each classification test, if RCBT was unable to complete the test in less than the cutoff time, it was terminated and it's run time was reported as the cutoff time. Hence, the 'RCBT' column gives lower bounds on RCBT's average run time per training set test. Finally, the '# RCBT DNF' column gives the number of tests RCBT was unable to finish in < the cutoff time, over the number of tests for which Top-K finished mining rule group upper bounds.

Explanation for varying nl values: Run time cutoffs were a necessity to mitigate excessive cross-validation CAR-mining times. Even with a cutoff of 2 hours these 100 PC experiments required

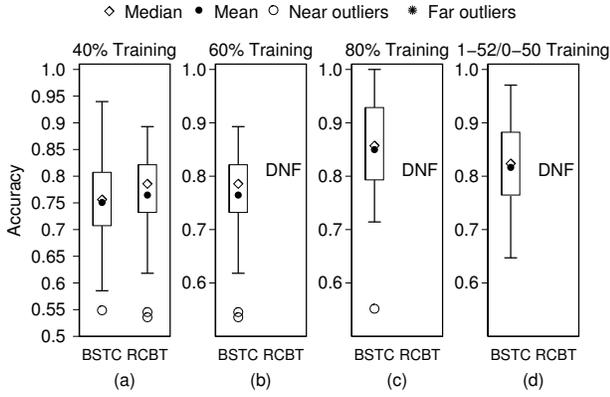


Figure 6: PC Cross-Validation Results

Training	BSTC	Top-k	RCBT	# RCBT DNF
40%	2.13	0.09	418.81	0/25
60%	4.93	5.06	≥ 7110.00	24/25
80%	5.78	120.63	≥ 7200 †	25/25 †
1-52/0-50	5.57	21.32	≥ 7200 †	25/25 †

Table 4: Average Run Times for the PC Tests (in seconds). Default cutoff time is 2 hours, and default nl value is 20. † indicates that the nl parameter was lowered to 2.

about 11 days of computation time. For the 80% and 1-52/0-50 training set sizes RCBT with $nl = 20$ failed to finish lower bound rule mining for all 50 tests within 2 hours. Thus, RCBT’s nl parameter was lowered from the default value of 20 to 2 in an attempt to improve its chances of completing tests. Not surprisingly, decreasing nl (i.e., mining fewer lower bound rules per Top-k rule group) decreases RCBT’s runtime. However, RCBT was still unable to finish lower bound rule mining for any tests.

Training	BSTC	RCBT
40%	75.08%	79.27%
60%	78.18%	85.45%
80%	84.98%	—
1-52/0-50	81.65%	—

Table 5: Mean Accuracies for the PC Tests that RCBT Finished. Both 40% and 60% results based on the tests that RCBT completed, namely 25 tests for 40% and 1 test for 60%.

Classification Accuracy: Figure 6 contains accuracy results for BSTC on all four Prostate Cancer test sets. Prostate Cancer box-plots for RCBT weren’t constructed for training set sizes RCBT was unable to complete all 25 tests for within the time cutoffs. In contrast, BSTC was able to complete each of the 100 PC classification tests in less than 6 seconds. Table 5 contains mean accuracies for the PC dataset with 40%, 60%, 80%, and 1-52/0-50 training. For each training set, the average accuracies were taken over the tests RCBT was able to complete within the cutoff time. Hence, the 40% row means were taken over all 25 results. Since RCBT was unable to complete any 80% or 1-52/0-50 training size tests we report these BSTC means over all 25 tests. RCBT has slightly better accuracy than BSTC on 40% training. For 60% training RCBT

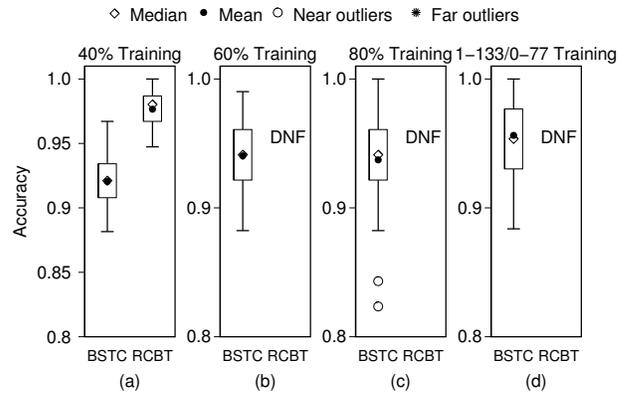


Figure 7: OC Cross-Validation Results

outperforms BSTC on the single test it could finish by more than 7%, although it should be kept in mind that RCBT’s results for the 24 unfinished tests could vary widely. Note that BSTC’s (mean) accuracy increases monotonically with training set size as expected. At 60% training BSTC’s accuracy behaves almost identically to RCBT’s 40% training accuracy (see Figure 6).

6.2.4 Ovarian Cancer (OC) Experiment

For the Ovarian Cancer dataset, which is the largest dataset in this collection, the Top-k mining method that is used by RCBT also runs into long computational times. Although Top-k is an exceptionally fast CAR group upper bound miner, it still depends on performing a pruned exponential search over the training sample subset space. Thus, as the number of training samples increases Top-k quickly becomes computationally challenging to tune/use.

Table 6 contains four average classification test run times (in seconds) for each Ovarian Cancer(OC) training size. As before, the second column run times each give the average time required to build both class 0/1 BSTs and then use them to classify all test’s samples with BSTC. Note that BSTC was able to complete each OC classification test in about 1 minute. In contrast, RCBT again failed to complete processing most classification tests within 2 hours.

Training	BSTC	Top-k	RCBT	# RCBT DNF
40%	30.89	0.6186	273.37	0/25
60%	61.28	41.21	≥ 5554.37	19/25
80%	71.84	≥ 1421.80	≥ 7205.43 †	21/22
1-133/0-77	70.38	≥ 1045.65	≥ 6362.86 †	20/23

Table 6: Average Run Times for the OC Tests (in seconds). Default cutoff time is 2 hours, and default RCBT nl value is 20. † indicates that the RCBT nl value was decreased to 2.

Table 6’s third column gives the average times required for Top-k to mine the top 10 covering rule groups upper bounds for each training set test (with the same 2 hour cutoff procedure as used for PC testing). The fourth column gives the average run times of RCBT on the tests for which Top-k finished mining rules (also with a 2 hour cutoff). Finally, the ‘# RCBT DNF’ column gives the number of tests that RCBT was unable to finish classifying in < 2 hours each, over the number of tests for which Top-k finished. Because RCBT couldn’t finish any 80% or 1-133/0-77 tests within 2 hours with $nl = 20$, we lowered nl to 2 for these training sizes.

Classification Accuracy: Figure 7 contains boxplots for BSTC on all four OC classification test sets. Boxplots were not generated for RCBT with 60%, 80%, or 1-133/0-77 training since it was unable to finish all 25 tests for all these training set sizes in < 2 hours each. Table 7 lists the mean accuracies of BSTC and RCBT over the tests on which RCBT was able to produce results. Hence, Table 7’s 40% row consists of averages over 25 results. Meanwhile Table 7’s 60% row results are from 6 tests, 80% contains a single test’s result, and 1-133/0-77 results from 3 tests. RCBT has better mean accuracy on the 40% training size, but the results are closer on the remaining sizes (< 4% difference over RCBT’s completed tests). Again, RCBT’s accuracy could vary widely on its uncompleted tests.

Training	BSTC	RCBT
40%	92.05%	97.66%
60%	95.75%	96.73%
80%	94.12%	98.04%
1-133/0-77	93.80%	96.12%

Table 7: Mean Accuracies for the OC Tests that RCBT Finished. Results based on the tests that RCBT completed, namely 25 tests for 40%, 6 tests for the 60%, 1 test for 80%, and 3 tests for 1-133/0-77.

CAR Mining Parameter Tuning and Scalability: We attempted to run Top-k to completion on the 3 OC 80% training and 2 OC 1-133/0-77 training tests which it couldn’t finish mining rules for within the 2 hour cutoff. Top-k finished two of the three 80% training tests in 775 min 43.64 sec (about 13 hours) and 185 min 3.29 sec. However, the third test ran for over 16,000 min (> 11 days) without finishing. Likewise, Top-k finished one of the two 1-133/0-77 tests in 126 min 45.15 sec but couldn’t finish the other in 16,000 min (> 11 days). After increasing Top-k’s support cutoff from 0.7 to 0.9 it was able to finish the two unfinished 80% and 1-133/0-77 training tests in 5 min 13.8 sec and 35 min 36.85 sec, respectively. However, RCBT (with $nl = 2$) then wasn’t able to finish lower bound rule mining for either of these two tests within 1,500 min. (more than a day). Clearly, *CAR-mining and parameter tuning on large training sets is computationally challenging*. As training set sizes increase, so will CAR-mining difficulties.

7. RELATED WORK

While operating on a microarray dataset, current CAR [9, 10, 30, 34] and other pattern/rule [20, 27] mining algorithms perform a pruned and/or compacted exponential search over either the space of gene subsets or the space of sample subsets. Hence, they are generally quite computationally expensive for datasets containing many training samples (or genes as the case may be). Part of the difficulty involved with mining CARs is that in addition to the exponentially large number of uninteresting rules that may be formed, there are usually many interesting rules as well. This means CAR miners such as CHARM [34] and CLOSET+ [30] may not only end up having to wade through a prohibitive number of low quality rules while discovering interesting CARs, but there may also be a huge number of repetitive CARs that are discovered.

The FARMER algorithm reduces the number of stored interesting rules by utilizing the notion of a **rule group**. Rule groups allow many interesting rules with similar sample support to be clustered together in a more compact form. Although rule groups provide a beneficial reduction in the number of interesting CARs which must be saved, there are typically still a large number of interesting rule

groups. Hence, for large datasets it can still be prohibitively expensive for FARMER to find and store all user targeted rule groups.

More recently, the Top-k algorithm has solved the problem of generating an excessive number of interesting (i.e. high confidence) user targeted rule groups. Top-k cleverly allows the user to decide on the number of best rule groups to find and store. Hence, a small number of non-redundant CAR rule groups may be stored and used for dataset analysis and classification. Although a significant step forward, Top-k still depends on performing a pruned exponential search of the dataset’s training sample subset space. Furthermore, the RCBT [9] classifier proposed by the Top-k authors requires a potentially prohibitively expensive breadth-first search on the subset space of antecedent genes in each discovered rule group upper bound.

BSTC is also related to decision tree-based classifiers such as random forest [6] and C4.5 family [26] methods. It is possible to represent any generalized boolean association rule as a decision tree and vice versa. However, it is generally unclear how the trees generated by current tree-based classifiers are related to high confidence/support CARs which are known to be particularly useful for microarray data [9, 10, 14, 18, 22]. BSTC is explicitly related to, and motivated by, CAR-based methods.

To the best of our knowledge there is no previous work on mining/classifying with BARs of the form we consider here. Perhaps the work closest to utilizing 100% BARs is the TOP-RULES [19] miner. TOP-RULES utilizes a data partitioning technique to compactly report item/gene subsets which are unique to each class set C_i . Hence, TOP-RULES discovers all 100% confident CARs in a dataset. However, the method must utilize an emerging pattern mining algorithm such as MBD-LLBORDER [13], and so generally isn’t polynomial time.

8. CONCLUSIONS AND FUTURE WORK

To address the computational difficulties involved with preclassification CAR mining (see Tables 4 and 6), we developed a novel method which considers a larger subset of CAR-related boolean association rules (BARs). These rules can be compactly captured in a table called Boolean Structure Table (BST), which can then be used to produce a BST classifier called BSTC. Comparison to the current best CAR classifier, RCBT, on several benchmark microarray datasets shows that BSTC is competitive with RCBT’s accuracy while avoiding the exponential costs incurred by CAR mining (see Section 6.2). Hence, *BSTC extends generalized CAR-based methods to larger datasets than previously practical*. Furthermore, unlike other association rule-based classifiers, *BSTC easily generalizes to multi-class gene expression datasets*.

BSTC’s per-query classification time: BSTC’s worst case theoretical per-query classification time is currently worse than a CAR-based method’s *after all exponential time CAR mining is completed* ($O(|S|^2 \cdot |G|)$ versus $O(|S| \cdot |G|)$). As future work we plan to investigate decreasing BSTC’s per-query classification time by carefully culling BST exclusion lists. For now we simply point out that BSTC’s Section 6 run times are reasonable and will remain so for larger problems on which CAR mining is infeasible (e.g., for OC training sets containing several hundred samples).

Generalizing BSTC: As future work we also plan to experiment with other boolean formula arithmetization procedures besides those employed to evaluate BST satisfaction levels in Algorithm 5. Multiple BST satisfaction level arithmetization procedures could be used along with a heuristic classification confidence measure employed to select the best one. One potential confidence measure is the normalized difference between the highest and second highest BST satisfaction level returned by each arithmetization

procedure. The larger the normalized difference, the more “sure” the procedure appears to be about its classification.

9. ACKNOWLEDGMENTS

We thank Anthony K.H. Tung and Xin Xu for sending us their discretized microarray data files and Top-k/RCBT executables. This research was supported in part by NSF grant DMS-0510203.

10. REFERENCES

- [1] The dprep package. <http://cran.r-project.org/doc/packages/dprep.pdf>.
- [2] R. Agrawal, T. Imielinski, and A. Swami. Mining associations between sets of items in large databases. *SIGMOD*, pages 207–216, 1993.
- [3] R. Agrawal and R. Srikant. Fast algorithms for mining association rules. *VLDB*, pages 487–499, 1994.
- [4] Z. Bar-Joseph. Analysing time series gene expression data. *Bioinformatics*, 20(16):2493–2503, 2004.
- [5] R. Bayardo and R. Agrawal. Mining the most interesting rules. *SIGKDD*, 1999.
- [6] L. Breiman. Random forests. *Mach. Learn.*, 45(1):5–32, 2001.
- [7] C. Chang and C. Lin. Libsvm: a library for support vector machines, 2001.
- [8] Y. Cheng and G. Church. Biclustering of expression data. *Proc. of the 8th Int. Conf. on Intelligent Systems for Molecular Biology (ISMB)*, 2000.
- [9] G. Cong, K. L. Tan, A. K. H. Tung, and X. Xu. Mining top-k covering rule groups for gene expression data. *SIGMOD*, 2005.
- [10] G. Cong, A. K. H. Tung, X. Xu, F. Pan, and J. Yang. Farmer: Finding interesting rule groups in microarray datasets. *SIGMOD*, 2004.
- [11] C. Cortes and V. Vapnik. Support-vector networks. *Machine Learning*, 20(3):273–297, 1995.
- [12] C. Creighton and S. Hanash. Mining gene expression databases for association rules. *Bioinformatics*, 19, 2003.
- [13] G. Dong and J. Li. Efficient mining of emerging patterns: discovering trends and differences. *KDD*, pages 43–52, 1999.
- [14] G. Dong, X. Zhang, L. Wong, and J. Li. Caep: Classification by aggregating emerging patterns. *Proc. 2nd Int. Conf. Discovery Science (DS)*, 1999.
- [15] J. Han, J. Pei, and Y. Yin. Mining frequent patterns without candidate generation. *SIGMOD*, 2000.
- [16] A. Icev, C. Ruiz, and E. F. Ryder. Distance-enhanced association rules for gene expression. *SIGKDD Workshop on Data Mining in Bioinformatics (BIOKDD)*, 2003.
- [17] J. Li, R. Topor, and H. Shen. Construct robust rule sets for classification. *KDD*, pages 564–569, 2002.
- [18] J. Li and L. Wong. Identifying good diagnostic genes or gene groups from gene expression data by using the concept of emerging patterns. *Bioinformatics*, 18:725–734, 2002.
- [19] J. Li, X. Zhang, G. Dong, K. Ramamohanarao, and Q. Sun. Efficient mining of high confidence association rules without support thresholds. *Principles of Data Mining and Knowledge Discovery (PKDD)*, pages 406 – 411, 1999.
- [20] W. Li, J. Han, and J. Pei. Cmar: Accurate and efficient classification based on multiple class-association rules. *ICDM*, 2001.
- [21] B. Liu, W. Hsu, and Y. Ma. Integrating classification and association rule mining. *KDD*, 1998.
- [22] T. McIntosh and S. Chawla. On discovery of maximal confident rules without support pruning in microarray data. *SIGKDD Workshop on Data Mining in Bioinformatics (BIOKDD)*, 2005.
- [23] R. Motwani and P. Raghavan. *Randomized Algorithms*. Cambridge University Press, 1995.
- [24] F. Pan, G. Cong, A. K. H. Tung, J. Yang, and M. J. Zaki. Carpenter: Finding closed patterns in long biological datasets. *KDD*, 2003.
- [25] J. Pei, X. Zhang, M. Cho, H. Wang, and P. S. Yu. Maple: A fast algorithm for maximal pattern-based clustering. *ICDM*, pages 259–266, 2003.
- [26] J. R. Quinlan. Bagging, boosting, and c4.5. *AAAI*, 1:725–730, 1996.
- [27] F. Rioult, J. F. Boulicaut, B. Cremilleux, and J. Besson. Using transposition for pattern discovery from microarray data. *DMKD*, pages 73–79, 2003.
- [28] R. Srikant and R. Agrawal. Mining quantitative association rules in large relational tables. *SIGMOD*, 1996.
- [29] S. Sudarsky. Fuzzy satisfiability. *International Conference on Industrial Fuzzy Control and Intelligent Systems (IFIS)*, 1993.
- [30] J. Wang, J. Han, and J. Pei. Closet+: Searching for the best strategies for mining frequent closed itemsets. *KDD*, 2003.
- [31] G. I. Webb. Discovering associations with numeric variables. *KDD*, pages 383–388, 2001.
- [32] G. I. Webb and S. Zhang. k-optimal-rule-discovery. *In Data Mining and Knowledge Discovery*, 10(1):39–79, 2005.
- [33] X. Xu, Y. Lu, A. K. Tung, and W. Wang. Mining shifting-and-scaling co-regulation patterns on gene expression profiles. *ICDE*, 2006.
- [34] M. Zaki and C. Hsiao. Charm: An efficient algorithm for closed association rule mining. *Proc. of the 2nd SIAM Int. Conf. on Data Mining (SDM)*, 2002.
- [35] H. Zhang, B. Padmanabhan, and A. Tuzhilin. On the discovery of significant statistical quantitative rules. *KDD*, pages 374–383, 2004.
- [36] L. Zhao and M. Zaki. Tricuster: An effective algorithm for mining coherent clusters in 3d microarray data. *SIGMOD*, pages 694–705, 2005.