

# Identifying Markov Blankets with Decision Tree Induction

## PrePublication Version

### **Lewis Frey**

Frey@vuse.vanderbilt.edu

Department of Biomedical Informatics, Vanderbilt University, 2209 Garland Avenue;  
Department of Electrical Engineering and Computer Science, Box 1679 Station B,  
Nashville, TN 37232 USA

### **Douglas Fisher**

DFisher@vuse.vanderbilt.edu

Department of Electrical Engineering and Computer Science, Vanderbilt University, Box  
1679 Station B, Nashville, TN 37235 USA

### **Ioannis Tsamardinos**

Ioannis.Tsamardinos@vanderbilt.edu

Department of Biomedical Informatics, Vanderbilt University, 2209 Garland Avenue,  
Nashville, TN 37232 USA

### **Constantin F. Aliferis**

Constantin.Aliferis@vanderbilt.edu

Department of Biomedical Informatics, Vanderbilt University

### **Alexander Statnikov**

Alexander.Statnikov@vanderbilt.edu

Department of Biomedical Informatics, Vanderbilt University

### **Abstract**

The Markov Blanket of a target variable is the minimum conditioning set of variables that makes the target independent of all other variables. Markov Blankets inform feature selection, aid in causal discovery and serve as a basis for scalable methods of constructing Bayesian networks. This paper applies decision tree induction to the task of Markov Blanket identification. Notably, we compare (a) C5.0, a widely used algorithm for decision rule induction, (b) C5C, which post-processes C5.0's rule set to retain the most frequently referenced variables and (c) PC, a standard method for Bayesian Network induction. C5C performs as well as or better than C5.0 and PC across a number of data sets. Our modest variation of an inexpensive, accurate, off-the-shelf induction engine mitigates the need for specialized procedures, and establishes baseline performance against which specialized algorithms can be compared.

Appears in the Proceedings of the Third IEEE International Conference on Data Mining, Melbourne, FL Nov 19-22 2003. pp 59-66

## 1. Introduction

The Markov Blanket (MB) is the minimum conditioning set that makes all other features independent for a particular target. Tsamardinos & Aliferis (2003) show that Markov Blankets consist of *strongly relevant* features as defined in relation to optimal classifiers (Kohavi & John, 1997). The Markov Blanket is particularly useful for data sets with large numbers of features. For example, gene expression data has tens of thousands of features. Arnone (1997) has estimated that on average, four to eight genes interact in higher organisms. Thus, an algorithm that finds the Markov Blanket of genes could reduce the data set by three orders of magnitude.

The Markov Blanket can be used for causal discovery under the condition of *faithfulness*, which is defined in Section 2. In essence, it can reduce the number of variables that need to be tested experimentally to discover direct causes of a target  $T$ .

Additionally, Markov Blanket discovery can be used for guiding Bayesian Network construction. The Markov Blanket for each variable is identified and used as a guide to construct the Bayesian Network for the domain. Margaritis and Thrun (1999) use this approach for Bayesian Network structure learning.

Several algorithms have been developed or proposed for identifying Markov Blankets (Tsamardinos, Aliferis & Statnikov 2003; Margaritis & Thrun, 1999; Koller & Sahami, 1996). Because of the relationship between Markov Blankets and feature relevance, and because decision tree induction has been used to perform feature selection (Cardie, 1993; Almaullim & Dietterich, 1991), we evaluate the ability of C5.0, an inexpensive, off-the-shelf induction engine, to identify Markov Blankets. Our evaluation motivates a modest post-processing step that demonstrably yields very accurate, low cost approximation of Markov Blankets. Our study thus forwards a scalable, accurate, and accessible algorithm for Markov Blanket identification. The observation that classification-motivated rule induction performs well at MB identification, also highlights the common principles of feature relevance that underlie induction of classifiers and Bayesian Networks (Tsamardinos & Aliferis, 2003).

The remainder of the paper reviews basic concepts of Bayesian Networks, describes C5.0, our variation called C5C, and the PC algorithm (Spirites, Glymour, & Clark, 2000), a prototypical algorithm for Bayesian Network construction, which we use for comparison purposes. An experimental comparison of PC, C5C, and C5.0 evaluates how accurately each method finds the Markov Blankets for target variables in a variety of domains. The paper concludes with a discussion of the results.

## 2. Bayesian Networks

A *Bayesian Network*,  $\langle V, G, J \rangle$ , consists of a set of variables  $V$ , a directed acyclic graph  $G$  and a joint probability distribution  $J$  over the variables  $V$  where the Markov Conditions holds (i.e., a variable is independent from all variables other than its descendants when conditioned on its parents) (Pearl, 1988).

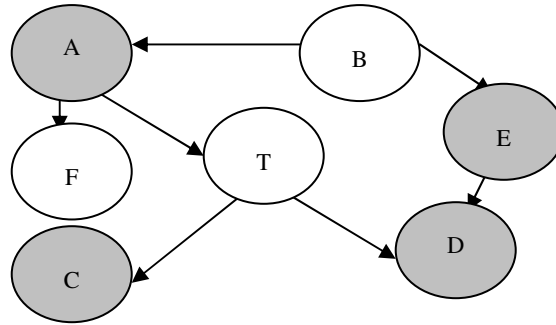


Figure 1. Example Bayesian Network Graph

The *Markov Blanket* (MB) of a target variable is the minimal set of variables  $X$  in  $V$  such that when a variable  $T$  in  $V$  is conditioned on  $X$ , any variable in  $F \subseteq V - X - T$  is independent of  $T$  (i.e.,  $P(T|XUF) = P(T|X)$ ).

In Bayesian Networks (BN) the union of parents and children of  $T$ , and parents of children (spouses) of  $T$  is equivalent to the Markov Blanket (Pearl, 1988). In Figure 1, the Markov Blanket for  $T$  is  $\{A, C, D, E\}$ . This means that variables  $B$  and  $F$  are independent of  $T$  conditioned on  $\{A, C, D, E\}$ .

A Bayesian Network  $N = \langle V, G, J \rangle$  is *faithful* to joint probability distribution  $J$  over feature set  $V$  if and only if all dependencies entailed by  $G$  and the Markov Condition are also present in  $J$ . A data-generating process  $K$  is faithful to  $N$ , if  $K$  in the sample limit produces data with joint probability distribution  $P$ , and  $N$  is faithful to  $P$ . A data set is faithful to a BN  $N$  if in the sample limit the data was generated by a data generation process that is faithful to  $N$ .

Some algorithms work well in identifying the MB from faithful data. If the data set is faithful to the BN, then the MB of each variable is unique (proof in appendix) and the comparison of the algorithms' ability to find the unique MB is valid. If the data is unfaithful, then the MB is not unique. This can occur, for example, when deterministic relationships exist in the data (discussion in Glymour and Cooper (1999)). Parity functions are examples of relationships that are not faithful. No current causal discovery methods can learn parity functions and most Bayesian methods assume faithfulness.

This paper focuses on the problem of inducing the Markov Blanket of the target concept in faithful distributions. For example, in the Bayesian Network depicted in Figure 1, the task of the algorithm is to predict  $A, C, D$  and  $E$  as the Markov Blanket and reject  $B$  and  $F$  as members of the Markov Blanket given a data set that is generated by the Bayesian Network. The algorithms for finding the Markov Blanket are discussed below.

### 3. Algorithms

This section describes the methods for Markov Blanket identification that we will compare. PC and C5.0 provide baseline performance to which C5C can be compared.

#### 3.1 PC

The PC algorithm (Spirtes, Glymour & Scheines, 2000) starts with a fully connected, unoriented Bayesian Network graph and goes through three phases. In phase I, the algorithm eliminates edges by using the criterion that variable  $A$  has a direct edge to variable  $B$  if, and only if, for all subsets of features there is no subset  $S$ , such that  $A$  is independent of  $B$  conditioned on  $S$ . In phases II and III, the algorithm orients the edges by performing global constraint propagation. If the algorithm is not able to orient some edges, the output is actually a class of structurally equivalent BNs. The PC algorithm uses significance thresholds based on the  $G^2$  statistic. In linear domains, Fisher's  $z$ -test is employed.

PC is intractable on densely connected data sets with a large number of variables. The algorithm has a complexity on the order of the number of variables raised to the maximal degree,  $d$ , for any node in the graph (i.e.,  $O(|V|^d)$ ). Consequently, it is exponentially bounded in the maximal degree. Similar complexity limitations hold for search and score Bayesian methods. The algorithm by Cheng et al. (2002) improves the complexity (to  $O(|V|^4)$ ) by introducing a distributional assumption ("monotone restriction") with properties currently being explored. In practice, these algorithms cannot run on more than a few hundred variables.

#### 3.2 C5.0

A C5.0 decision tree is constructed using *GainRatio*. *GainRatio* is a measure incorporating entropy. Entropy ( $E(S)$ , Eq. 1) measures how unordered the data set is. It is denoted by the following equation when there are classes  $C_1, \dots, C_N$  in data set  $S$  where  $P(S_c)$  is the probability of class  $C$  occurring in the data set  $S$ :

$$E(S) = - \sum_{c=1}^N P(S_c) * \log_2 P(S_c) \quad (1)$$

*Information Gain* is a measure of the improvement in the amount of order.

$$Gain(S, V) = E(S) - \sum_{v \in \text{values}(V)} \frac{|S_v|}{|S|} * E(S_v) \quad (2)$$

*Gain* has a bias towards variables with many values that partition the data set into smaller ordered sets. In order to reduce this bias, the entropy of each variable over its  $m$  variable values is calculated as *SplitInfo* (Eq. 3).

$$SplitInfo(S, V) = \sum_{i=1}^m - \frac{|S_i|}{|S|} * \log_2 \frac{|S_i|}{|S|} \quad (3)$$

*GainRatio* (Eq. 4) is calculated by dividing *Gain* by *SplitInfo* so that the bias towards variables with large value sets is dampened.

$$GainRatio(S, V) = \frac{Gain(S, V)}{SplitInfo(S, V)}$$

C5.0 builds a decision tree greedily by splitting the data on the variable that maximizes gain ratio. CF-pruning is set to 25% (the default setting for C5.0) in this paper.

A final decision tree is changed to a set of rules by converting the paths into conjunctive rules and pruning them to improve classification accuracy. When using C5.0 to predict the Markov Blanket, each variable involved in any rule output by C5.0 is predicted to be in the Markov Blanket of the target/class variable.

### 3.3 C5C

The working hypothesis is that frequently occurring features in C5.0 production rules provide a good approximation of the  $MB(T)$ . This hypothesis is encapsulated as a simple augmentation of C5.0 to count the number of variable occurrences in C5.0 rules and use the most frequent to predict the Markov Blanket. This is called the C5C algorithm and it uses the C5.0 rules output to order the variables from most likely MB variables to least likely.. The modification consists of a simple script that counts the occurrence of the variables in the C5.0 rules output (i.e., if variable occurs in rule, then increment the variable's frequency count) A frequency threshold is chosen to distinguish between Markov and non-Markov Blanket variables.

Three methods will be used to identify the MB given the C5C ranking of variables. First is choosing the best frequency threshold given knowledge of the true MB (this is used as a comparison in sections 6.1 and 6.2 and is intended as a best-case analysis). The second strategy finds the frequency threshold that gives the best accuracy of C5.0 on a test set. The third strategy employs the  $G^2$ -test to test for independence.

## 4. Experimental Methods

The experiments involve generating data from six Bayesian Networks and comparing the prediction of Markov Blankets by C5.0, C5C and PC. The first five data sets are publicly available Bayesian Networks: Alarm (Beinlich et al., 1989), Hailfinder (Abramsom et al., 1996), Insurance (Binder et al., 1997), Mildew and Barley (Kristensen & Rasmussen, 2002). These data sets are generated using a logic sampling data-generating approach, implemented in the Hugin package (Andersen et al, 1989), run to produce a data set of 20,000 discrete valued instances. The last data set is an Artificial Bayesian Network that is constructed to examine the effect of sample size on learning.

Using these generated data sets, C5.0, C5C and PC's abilities to reconstruct the Markov Blanket for each variable in the data sets are compared. Since the original networks are available, the predicted Markov Blankets can be compared against the actual Markov Blankets in terms of sensitivity and specificity. *Sensitivity* is the ratio of correctly predicted MB variables over true MB variables. *Specificity* is the ratio of correctly predicted (i.e., excluded) non-MB variables over true non-MB variables. Ratios for predicted Markov Blankets of different sizes and characteristics are compared.

In comparing C5.0 and C5C, a measure of closeness to the true MB is determined by calculating the distance (*dist*) of sensitivity (*sen*) and specificity (*spec*) of a predicted MB to that of the true MB (Eq. 5). Thus, the smaller the distance, the closer the predicted MB is to the true MB.

$$dist = \sqrt{(1 - sen)^2 + (1 - spec)^2} \quad (5)$$

The area under the Receiver Operating Characteristic (ROC) curve is calculated and serves as a measure of how well C5C and PC output the Markov Blanket with different thresholds. For C5C, we vary the frequency threshold, and for PC we vary the significance threshold. A one hundred percent area under the ROC curve equates to exactly identifying the Markov Blanket and fifty percent area is guessing.

## **5. Data Sets**

### **5.1 Alarm Network**

A data set generated from the Alarm Network, a medical monitoring network, is used to explore how well the algorithms find the Markov Blankets for all of the variables in the network. The alarm network has 37 variables, each of which has its own Markov Blanket.

### **5.2 Hailfinder Network**

The data set generated from the Hailfinder Network has 56 variables, which are used for severe weather forecasting (Abramsom et al., 1996).

### **5.3 Insurance Network**

The Insurance network has 27 variables and is used for estimating expected claim costs for insurance policies (Binder et al., 1997).

### **5.4 Mildew Network**

The data set generated from the Mildew Network has 35 variables, which determine amounts of fungicides to use against attacks of mildew on wheat.

### **5.5 Barley Network**

The Barley data has 48 variables, which help predict yield and quality parameters for growing barley without pesticides (Kristensen & Rasmussen, 2002).

### **5.6 Artificial Bayesian Network**

A data set generated from an artificial Bayesian Network is used to explore the effect of the number of variables and the sample size upon the algorithms' ability to find the Markov Blanket. For this data set, the Markov Blanket is found for only one variable. The Markov Blanket for this variable has three parents, two children and one parent of a child for a total of six variables. The Markov Blanket for this variable is kept constant over the conditions of increasing noisy-variables and reduced sample size.

The above data sets are examined because their true Markov Blankets are known. Because the predicted Markov Blankets of the methods are compared against the true Markov Blankets, this is necessary.

## **6. Results**

### **6.1 Comparisons between C5C and C5.0**

In Table 1, the average sensitivity, specificity and distance over target variables for C5.0 and the best thresholds for C5C are listed for the first five Bayesian Networks. Table 1

illustrates that C5C's average predicted MB is closest to the true MB. Examining the sensitivity and specificity values, it appears that C5C's distances are improved by the increase in specificity without the loss of sensitivity. Discarding low frequency variables from the MB usually causes non-MB variables to be discarded.

*Table 1.* Average over target variables of sensitivity (sen), specificity (spec) and distance (dist) for C5.0 and the best thresholds for C5C. Data sets have 20,000 instances. The asterisk (\*) denotes the mean distance for C5C is significantly different from C5.0 by the paired Wilcoxon signed rank test of the equality of means ( $p < 0.05$ ).

DATA SET	C5.0			C5C		
	Sen	Spec	Dist	Sen	Spec	Dist
ALARM	0.83	0.84	0.32	0.82	0.99	0.18*
HAILFINDER	0.81	0.27	0.89	0.80	0.98	0.22*
INSURANCE	0.89	0.46	0.64	0.78	0.93	0.25*
MILDEW	0.94	0.11	0.95	0.80	0.90	0.26*
BARLEY	0.84	0.43	0.67	0.76	0.94	0.28*

Table 2 shows that for 156 variables (75% of the total variables) C5C's predicted MB is closer to the true MB than C5.0's predicted MB. C5C is equal to C5.0 for the remaining 47 variables.

*Table 2.* Number of target variables (var) out of the total for each data set that the distance of C5C's best predicted MB is closer to the true MB than C5.0.

DATA SET	FREQ THAT C5C IS CLOSER TO TRUE MB THAN C5.0	TOTAL VAR
ALARM	15	37
HAILFINDER	46	56
INSURANCE	22	27
MILDEW	34	35
BARLEY	39	48

Table 3 shows the improvement in average distance over the 156 variables where C5C is closer to the true MB than C5.0. It also shows the reduction in the average predicted MB size from C5.0 to the best C5C thresholds.

*Table 3.* Average Markov Blanket size and average distance to the true MB for C5.0 and best threshold for C5C. The asterisk (\*) denotes significance by the paired Wilcoxon signed rank test of the equality of means ( $p < 0.05$ ) in comparing C5.0 and C5C distance.

DATA SET	AVG MB SIZE		DISTANCE	
	C5.0	C5C	C5.0	C5C
ALARM	16	3	0.40	0.05*
HAILFINDER	49	4	0.93	0.12*
INSURANCE	19	6	0.67	0.20*
MILDEW	31	6	0.98	0.27*
BARLEY	34	7	0.73	0.22*

Table 3 also demonstrates that C5C reduces the size of the predicted MB more than C5.0. The quality of the predicted MB is of interest. This also supports the hypothesis

that the low frequency variables tend to be non-Markov Blanket variables. This is why there can be such a reduction in the predicted MB size while improving the average distance.

## **6.2 C5C Identifying Markov Blanket variables**

The results in Tables 1-3 show that C5C tends to rank MB variables higher than non-MB variables. Two methods that use C5C ranking of variables are tested in their ability to identify the MB: method one uses decision tree test set accuracy and method two uses the  $G^2$  test.

For method one, C5.0 decision tree test set accuracy is assumed to improve when only MB variables are used in the training set. The filtered variable set with the best accuracy should correspond to the Markov Blanket because the Markov Blanket is the complete set of strongly relevant features. Variables are added to the training set for the decision tree starting with the variable(s) with the highest C5C ranking and proceeding to the lowest via a wrapper method (Caruana & Freitag, 1994). The accuracy of the classifier is recorded for each of these subsets of variables. The subset with the highest test accuracy is proposed as the MB. The results in Table 4 are for running C5.0 with training set with sample size 16,000 for the training set and 4,000 for the test set.

The second method uses the  $G^2$ -test (Spirites, Glymour & Scheines, 2000), to test for independence between two variables given a conditioning set ( $p < 0.05$ ). A weakness of the  $G^2$ -test is that its validity is sensitive to the amount of sample. The larger the Markov Blanket conditioning set the more sample needed (i.e., sample grows exponentially with conditioning set size). For this method all variables with zero C5C frequency counts are classified as non-Markov Blanket variables. Then the  $n$  top ranked variables of C5C are proposed as the conditioning set for the  $G^2$ -test. The top  $k$  ranked variables (where  $k > n$ ) are then tested for independence from the target variable given the conditioning set. This is an heuristic approach to identify a larger MB than the sample size supports for the  $G^2$ -test. If the variables are not independent, they are included in the final MB. Each variable in the conditioning set is then tested for independence given the other features in the conditioning set. Those that are not independent are included in the final MB. In Table 4  $n=4$  and  $k=10$  and the sample size is 20,000. The conditioning set is size 4 because with a sample of 20,000 the  $G^2$ -test tends to be valid for the data sets examined. This can be increased with larger sample size or reduced with smaller sample size.

One of the benefits of the ordered C5C output is that it provides a subset of the variables (i.e.,  $k$  variables with high frequency counts) to be examined for Markov Blanket candidacy. For the examined data sets, the parameter  $k$  is arbitrarily set to 10. To minimize the exclusion of Markov Blanket variables, it is recommended that  $k$  is greater than the expected quantity of Markov Blanket variables. For example, in some biological data set, Arnone (1997), there are on average less than eight interacting variables, so  $k = 10$  is reasonable and  $k = 15$  would be more conservative. Since the assumption of C5C is that variables in the C5.0 rules reflect the Markov Blanket, all variables with zero frequency counts (i.e., they do not occur in C5.0 rules) are not considered candidates for the Markov Blanket. For some data sets examined, many of the variables in the C5C output have zero frequency.

Table 4. Average over target variables of sensitivity (sen), specificity (spec) and distance (dist) for C5C with the C5.0 decision tree test set accuracy determining threshold (left) and the  $G^2$ -test identifying the MB (right). For the test accuracy the training set is 16,000 instances and the test set is 4,000 instances. The  $G^2$ -test uses 20,000 instances. The asterisk (\*) denotes the mean distance for the method is significantly different from the mean distance of C5.0 (Table 1) by the paired Wilcoxon signed rank test of the equality of means ( $p < 0.05$ ) the plus (+) is marginal at  $p < 0.1$ .

DATA SET	C5C – TEST ACC.			C5C – $G^2$		
	Sen	Spec	Dist	Sen	Spec	Dist
ALARM	0.79	0.97	0.23	0.83	0.97	0.20 <sup>+</sup>
HAILFINDER	0.76	0.95	0.27*	0.79	0.88	0.29*
INSURANCE	0.74	0.80	0.42*	0.76	0.88	0.31*
MILDEW	0.75	0.85	0.38*	0.78	0.80	0.35*
BARLEY	0.71	0.91	0.36*	0.76	0.87	0.31*

Table 4 shows that the two methods improve the identification of the MB over the C5.0 algorithm (see Table 1 C5.0 dist column). For some data sets the distances are close to the “best distance” (Table 1 C5C dist column), but both methods are significantly different from the “best distance”. In comparing the two methods against each other they are not significantly different ( $p < 0.05$ ) except for Hailfinder and Insurance. This means the two methods are comparable in their ability to identify the MB.

### 6.3 Comparisons between C5C and PC

C5C and PC are run on the Bayesian Network data for each variable. The average area under the ROC is then computed for the algorithms by plotting the sensitivity and specificity values for the thresholds.

Six thresholds are used with C5C. The first threshold proposes all variables in the C5.0 rules as the Markov Blanket for the target variable (0% of variables in rules are excluded). The second threshold excludes from the Markov Blanket variables in the C5.0 rules whose counts are infrequent (20% of the most frequent feature count for the given target variable). For example, if the most frequent variable in the rules occurred 100 times, then the second threshold would exclude any variable in the C5.0 rules that occurred fewer than 20 times. The thresholds occur in 20% increments (i.e., 20%, 40%, 60%, 80% and 100%). The area under the ROC curve drops if C5.0 rules don’t contain Markov Blanket variables or if infrequent variables in the C5.0 rules are part of the Markov Blanket.

It is computationally cheap to get multiple thresholds from C5.0’s output file, unlike PC that must be run for each threshold, so a finer granularity of thresholds (one hundred and one) are examined to see if better ROC area can be obtained. The thresholds are incremented by 1% of the maximum count variable in the C5.0 rules.

The six thresholds for the PC algorithm are 0.005, 0.01, 0.02, 0.05, 0.10 and 0.20. These thresholds are significance levels for the  $G^2$  test used by the algorithm. The smaller the threshold, the higher the confidence level that the variables are independent. This means the smaller the threshold, the fewer the number of edges removed from the fully

connected graph used by PC to determine Markov Blankets. PC must be run separately for each threshold confidence level.

Table 5 shows that both PC and C5C performed comparatively well in finding the Markov Blankets for Alarm, Hailfinder and Insurance data sets (i.e., not significantly different) across both sets of thresholds.

C5C does better than PC on both the Mildew and Barley data sets for both thresholds (Table 5). This is because these data sets have variables with large value sets (100 values). This condition makes it difficult for PC to eliminate edges from the fully connected graph.

*Table 5. Average ROC over all variables in Network: PC & C5C with 6 thresholds (T'holds) and C5C with 101 thresholds. Data sets have 20,000 instances. The asterisk (\*) denotes the mean ROC for C5C as significantly different from PC by the paired Wilcoxon signed rank test of the equality of means ( $p < 0.05$ ).*

DATA SET	TOTAL TARGETS	PC	C5C	
			6	101
# OF T'HOLDS		6	6	101
ALARM	37	96.3	91.1	91.4
HAILFINDER	56	91.7	87.5	88.9
INSURANCE	27	82.0	86.3	88.0
MILDEW	35	64.3	80.0*	88.0*
BARLEY	48	50.0	81.6*	85.9*

In summary, Table 5 shows that C5C is not significantly different from PC for three data sets and is significantly better for two data sets. C5C behaves similarly with 6 and 101 thresholds so the distinction will be dropped in further analysis. Importantly, C5C can be run on data with many more variables than is reasonable with PC.

#### **6.4 C5C over increasing noisy-variables and reduced sample size**

Figure 2 shows that PC and C5C have similar performance across sample size and number of variables. C5C appears to do better at smaller sample size across the number of variables. C5C's greedy nature works in its favor while PC tends to get lower ROC for small sample size because it starts with a fully connected graph and needs enough evidence to remove edges. If there is not enough evidence to remove edges, then specificity is low.

When there are a larger number of noisy variables (e.g., 1,000 variables), the C5C algorithm needs more data to distinguish the MB. Still, as Figure 2 demonstrates, C5C is able to find the Markov blanket for large sample and larger variable size.

C5C is also time-efficient. The largest computation (i.e., 1,000 variables for a sample size of 20,000) takes 30 seconds to perform on a 2.4GHz Xeon desktop machine for all thresholds.

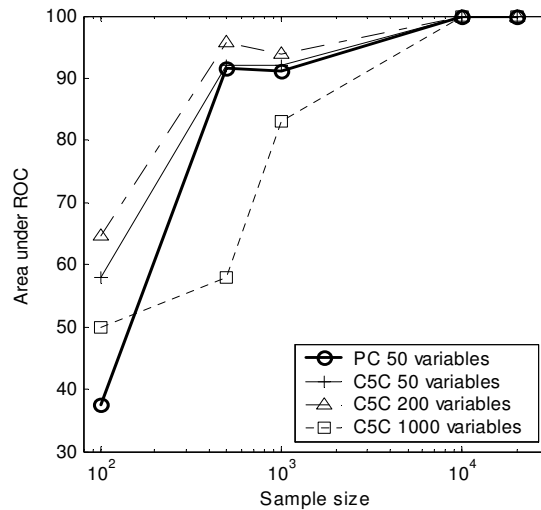


Figure 2. Area under the ROC for MB size 6 for PC with 50 variables and C5C with 50, 200 and 1,000 variables. The artificial Bayesian Network sample sizes are 100, 500, 1,000, 10,000 and 20,000 examples.

## 7. Discussion

C5C is a simple modification to C5.0 that approximates the Markov Blanket as well as or better than C5.0 and PC in these experiments. Additionally, for the constructed Bayesian Network there is evidence that smaller sample sizes have good values of area under the ROC curve. This reduction in sample size could be due to the greedy nature of the algorithm. It chooses the best variables (usually Markov Blanket variables) given the small sample. The C5C algorithm is quick, in that it greedily hill climbs through the space of possible decision trees.

Note that the concept of causal structure within a data set is, in general, different than the goal of improving classification accuracy. C5.0 is designed to improve classification accuracy. Thus, there is no explicit bias encoded into the algorithm for finding causal structure. C5.0 outputs a decision tree and a set of rules with the goal of learning to classify the data set with highest accuracy.

### 7.1 Limitations

Due to this bias, the C5C algorithm is not able to find the Markov Blanket for target variables that occur a disproportionate number of times in one class. In such a case, C5.0 implementation uses default classification and assigns the variable to the majority class. This happens once in Insurance, twice in Barley and three times in Alarm. When a default classification occurs, no rules are formed so no variables can be part of the Markov Blanket. This cannot be remedied by changing the pruning level.

The C5C algorithm is not able to predict the Markov Blanket when one variable predicts the target variable without error, even though there are other variables in the Markov Blanket. The single variable becomes the only variable in the rules so all other variables are excluded from the Markov Blanket. Changing the pruning level does not affect the rules when one variable perfectly predicts the target variable. This happens for

eight variables in Hailfinder. However, this is an unfaithful distribution because there are deterministic relationships in the data set.

A related approach is that of Li et al. (2001) who use Genetic Algorithms and frequently occurring features in their classifiers to determine feature relevance in gene expression data. Their motivation is establishing a set of predictors that are robust to noise. C5C, even in non-noisy data, gives a good approximation of the MB.

## 7.2 Future Work

Recall that the task of feature selection is to find a subset that maintains the optimal classifier. Markov blankets have been mapped to the concept of strongly relevant features in feature selection (Tsamardinos & Aliferis, 2003). Hence, this simple and scalable method for finding the Markov Blanket is widely applicable in the area of feature selection in machine learning and data mining.

The method of using the C5C ranking of features to identify the minimal feature subset for improvement in classification accuracy had some success (Table 4). The thresholds of C5C produce different filtered feature sets. These feature sets can be used as input into the C5.0 decision tree algorithm with cross-validated accuracy. The filtered feature set with the best accuracy should correspond to the Markov Blanket because the Markov Blanket is the complete set of strongly relevant features.  $G^2$ -test method generally performed better than the decision tree accuracy test. The MB from this test could also be used as a filtered feature set to improve classification accuracy.

Additional areas of future exploration are methods for selecting pruning levels and weighting the importance of features as opposed to straight counting. The weighting could include factors such as size of rule, accuracy of rule and its coverage.

## 8. Conclusion

It has been shown empirically that a simple post-processing step using the C5C algorithm performs at least as well as the PC algorithm when finding the Markov Blanket on the data sets compared. The method of using the ranked C5C variables in conjunction with performing  $G^2$ -tests (i.e., the cheapest of the two methods that we examined) can be used to identify the Markov Blanket variables. In addition, C5.0 scales up very well. Since C5.0 can be used to obtain the Markov Blanket right out of the box, applying this counting algorithm, C5C, for useful applications is a simple matter.

## Acknowledgements

Support for this research was provided in part by NIH grant LM 007613-01.

## Appendix

Denote d-separation of  $X$  and  $Y$  by  $Z$  in the graph of BN  $C$ , as  $D_C(X; Y | Z)$ . Denote independent  $X$  and  $Y$  given  $Z$  in the probability distribution as  $I(X; Y | Z)$ . Denote the set of parents and children of  $T$  in the graph of BN  $C$  as  $PC_C(T)$ . Denote the set of the parents, children, and spouses of  $T$  in the graph of BN  $C$  as  $MB_C(T)$ .

Proposition 1: A faithful BN  $C$  is a perfect map of independencies and dependencies of the data. In other words, the d-separation criterion gives all the dependencies and independencies:  $D_C(X; Y | Z) \Leftrightarrow I(X; Y | Z)$  (Pearl, 1988).

Proposition 2: In a faithful BN  $C$  on variables  $V$ ,  $X$  and  $T$  have a direct edge between them, if and only if  $\forall S \subseteq V, \neg I(X;T|S)$  (Spirtes et al., 2000).

Proposition 3: In any BN  $C$  on variables  $V$ , the set of parents, children, and spouses of  $T$  d-separates  $T$  from any other node in  $V$  (Neapolitan, 1990).

Theorem 1. For any two BNs  $C$  and  $N$ , both faithful to the same distribution  $J$  on a variable set  $V$ , and for any variable  $T$  in that distribution,  $PC_C(T) = PC_N(T)$  (i.e., the set  $PC(T)$  is unique in all faithful BNs).

Proof: A variable  $X$  has an edge to  $T$  in the graph of  $C$ , if and only if, for every subset  $S$ , it is the case that  $\neg I(X;T|S)$  (prop 2). Suppose that  $X \in PC_C(T)$ , but  $X \notin PC_N(T)$ .  $X \in PC_C(T)$  implies that  $\forall S \subseteq V, \neg I(X;T|S)$ .  $X \notin PC_N(T)$  implies that  $\exists S \subseteq V, I(X;T|S)$ , which contradicts the previous statement. <sup>1</sup>

So the index  $PC_C(T)$  can be replaced by  $PC(T)$  when referring to a faithful distribution network  $C$ .

Theorem 2: For any two BNs  $C$  and  $N$ , both faithful to the same distribution  $J$  on a variable set  $V$ , and for any variable  $T$  in that distribution,  $MB_C(T) = MB_N(T)$  (i.e., the set  $MB(T)$  is unique for all faithful BNs).

Proof: By definition  $MB_C(T) = PC(T) \cup Spouses_C(T)$  and  $MB_N(T) = PC(T) \cup Spouses_N(T)$ , where  $Spouses_C(T)$  are the spouses of  $T$  in the graph of  $C$ .

Let  $X \in MB_C(T)$ . The following shows that  $X \in MB_N(T)$ . If  $X \in PC_C(T)$ , by Theorem 1 then  $X \in PC_N(T)$  and so  $X \in MB_N(T)$ . Thus, only the case where  $X \in Spouses_C(T) \setminus PC_C(T)$  needs to be considered.

Let  $X \in Spouses_C(T) \setminus PC_C(T)$ , i.e., it is a spouse but not a parent or a child, and so there is no direct edge between  $X$  and  $T$ . Suppose  $X \notin MB_N(T)$ , then this leads to contradiction. Since  $X \notin MB_N(T)$ , then  $D_N(X;T|MB_N(T))$  because the  $MB_N(T)$  d-separates every other node from  $T$  (prop 3). Thus,  $I(X;T|MB_N(T))$  (prop 1).

Now, there is a subset  $S \subseteq V$  such that  $I(X;T|S)$  (prop 2) since  $X$  and  $T$  do not have a direct edge. Since  $X$  is a spouse, there is a least one common child of  $X$  and  $T$  in the graph of  $C$ . Let this be called  $Y$ . Conditioned on  $Y$  and *any* superset of  $Y$ ,  $X$  and  $T$  have to be d-connected (i.e., not d-separated) and so they also have to be dependent, i.e.,  $\forall S \subseteq V - \{X\}, \neg I(X, T|S \cup \{Y\})$  (by prop 1 and definition of d-separation). However, by theorem 1,  $Y$  is a member of  $PC(T)$  and so a member of  $MB_N(T)$ . Thus, it holds that  $\neg I(X, T|MB_N(T))$  (since  $MB_N(T)$  a superset of  $Y$  that does not contain  $X$ ). This contradicts are statement above that  $I(X;T|MB_N(T))$ .

With the symmetric argument if  $X \in MB_N(T)$ , then it has to be  $X \in MB_C(T)$  and so  $MB_C(T) = MB_N(T)$  <sup>1</sup>

## References

- Abramson, B., Brown, J., Edwards, W., Murphy, A. and Winkler, R.L. (1996).  
Hailfinder: A Bayesian system for forecasting severe weather. *International Journal of Forecasting* , 12, 57-71.

- Aluallim, H., and Dietterich, T. G. (1991). Learning with many irrelevant features. *Proceedings, Ninth National Conference on Artificial Intelligence*, pp. 547-552. Anaheim, CA. AAAI Press/ The MIT Press.
- Arnone, M.I. and Davidson, E.H. (1997). The hardwiring of development: organization and function of genomic regulatory systems. *Development*, 12(4), pp. 1851-1864.
- Andersen, S.K., Olesen, K.G., Jensen, F. V. and Jensen, F. (1989). HUGIN - a shell for building bayesian belief universes for expert systems in *Proceedings of the Eleventh International Joint Conference on Artificial Intelligence* (pp. 1080-1085).
- Beinlich, I., Suermondt, G., Chavez R. and Cooper G. (1989). The ALARM monitoring system: A case study with two probabilistic inference techniques for belief networks, *Proceedings of 2'nd European Conference on AI and Medicine*, Springer-Verlag, Berlin.
- Binder, J., Koller, D., Russell, S. and Kanazawa, K. (1997). Adaptive Probabilistic Networks with Hidden Variables. *Machine Learning*, 29, 213-244.
- Cardie, C. (1993). Using decision trees to improve case-based learning, in *Proceedings of the Tenth International Conference on Machine Learning* (pp. 25-32). Morgan Kaufmann.
- Caruana, R. and D. Freitag (1994). Greedy Attribute Selection, in *International Conference on Machine Learning*.
- Cheng, J., Greiner, R., Kelly, J., Bell D. and Liu, W. (2002). Learning Bayesian networks from data: an information-theory based approach. *Artificial Intelligence 137*, pp. 43-90.
- Glymour, C. & Cooper, G.F. (1999). *Computation, Causation and Discovery*. AAAI/The MIT Press.
- Kohavi, R. and G.H. John (1997). Wrappers for Feature Subset Selection. *Artificial Intelligence*, 97(1-2), 273-324.
- Kristensen, K. and Rasmussen, I.A. (2002). The use of a Bayesian network in the design of a decision support system for growing malting barley without use of pesticides. *Computers and Electronics in Agriculture*, 33, 197-217.
- Koller, D and M. Sahami (1996). Toward Optimal Feature Selection in *Thirteenth International Conference in Machine Learning* (pp. 284-292).
- Li, L., Pedersen, L.G., Darden, T.A. and Weinberg, C.R. (2001). Computational Analysis of Leukemia Micorarray Expression Data Using the GA/KNN Method. CAMDA'01.
- Margaritis, D. and Thrun, S. (1999). Bayesian Network Induction via Local Neighborhoods. Technical Report: CMU-CS-99-134.
- Neapolitan, R.E. (1990). *Probabilistic Reasoning in Expert Systems: Theory and Algorithms*. John Wiley and Sons.
- Pearl, J. (1988). *Probabilistic Reasoning in Intelligent Systems*. San Mateo, CA: Morgan Kaufman.

- Spirtes, P., C. Glymour, and R. Scheines (2000). *Causation, Prediction and Search*. Cambridge, MA: MIT Press.
- Quinlan, J.R. (1987). Induction of decision trees. *Machine Learning* 1 pp. 81-106.
- Tsamardinos, I. and Aliferis, C.F. (2003) Towards Principled Feature Selection: Relevancy, Filters and Wrappers in *Proceedings of the Ninth International Workshop on Artificial Intelligence and Statistics*.
- Tsamardinos, I., Aliferis, C.F. and Statnikov, A. (2003) Algorithms for Large Scale Markov Blanket Discovery to appear in *Proceedings of the 16th International FLAIRS Conference*.