

Improving the K2 Algorithm Using Association Rule Parameters

Evelina Lamma

Fabrizio Riguzzi

Sergio Storari

University of Ferrara, Department of Engineering,
via Saragat 1, 44100 Ferrara, Italy
{elamma,friguzzi,sstorari}@ing.unife.it

Abstract

A Bayesian network is an appropriate tool to work with a sort of uncertainty and probability, that are typical of real-life applications. Bayesian network arcs represent statistical dependence between different variables. In the data mining field, association rules can be interpreted as well as expressing statistical dependence relations. K2 is a well-known algorithm which is able to learn Bayesian network. In this paper we present an extension of K2 called K2-rules that exploits a parameter normally defined in relation to association rules for learning Bayesian networks. We compare K2-rules with K2 and SLA on the problems of learning four Bayesian networks. The experiments show that K2-rules improves both K2 and SLA with respect to the quality of the learned network and K2 with respect to the execution time.

Keywords: Bayesian Networks, Machine Learning, Association Rules

1 Introduction

A Bayesian network [19],[14] is an appropriate tool to work with uncertainty and probability, that are typical of real-life applications. A Bayesian network is a directed, acyclic graph (DAG) whose nodes represent random variables. In a Bayesian network each node is conditionally independent from any subset of nodes that are not its descendants, given its parents.

By means of Bayesian networks, we can use information about the values of some variables

to obtain probabilities for other variables. A probabilistic inference takes place once the probabilities functions of each node conditioned to just its parents are given. These are usually represented in a tabled form, named Conditional Probability Table (CPT).

Given a training set of examples, learning a Bayesian network is the problem of finding the structure of the direct acyclic graph and the CPT associated with each node that best match (according to some scoring metric) the dataset. Optimality is evaluated with respect to a given scoring metric (e.g., description length or posterior probability [3],[9],[11],[12],[13],[17],[20],[25]). A procedure for searching among possible structures is needed. However, the search space is so vast that any kind of exhaustive search cannot be considered, and often a greedy approach is followed.

The K2 algorithm [9] is a typical search and score method. It starts by assuming that a node has no parents, after which, in every step it adds incrementally the parent whose addition mostly increases the probability of the resulting structure. K2 stops adding parents to the nodes when the addition of a single parent cannot increase the probability of the network given the data. Other search and score methods include the MDL algorithm [25], and the CB algorithm [22].

In this work, we propose the algorithm K2-rules in order to improve the quality of learned networks and to reduce the computational resources needed. This algorithm uses data mining techniques, and in particular the computation of parameters normally defined in relation to association rules [1], to obtain new knowledge to be used for improving some of the steps of K2. Association rules describe

correlation of events, and can be viewed as probabilistic rules. Two events are “correlated” when they are frequently observed together. Both Bayesian network arcs and association rules represent dependence relations among variables so it may be interesting to integrate these methodologies in order to improve Bayesian network learning. Each association rule is characterized by several parameters which can be used to identify the absence of dependence among the nodes. In this work, we exploit in particular the leverage parameter in order to improve K2.

The paper presents the results of a comparison between K2, K2-rules and SLA [7], another well-known learning algorithm. SLA computes the mutual information of each pair of nodes as a measure of dependency and creates the network using this information.

The paper is structured as follows. Section 2 describes the K2 algorithm. In Section 3, we briefly present association rules. Section 4 illustrates the algorithm K2-rules. In Section 5, we show an experimental comparison between K2, K2-rules and SLA considering four of the most known Bayesian networks. Section 6 discusses related works. Finally, in Section 7, we conclude and present future work.

2 The K2 algorithm

In literature, we find different approaches for Bayesian network learning. Some of them are based on the search and score methodology [3],[9],[11],[12],[13],[17],[20],[25], and the others follow an information theory based approach [7],[22].

A procedure frequently used for learning the structure of a Bayesian network from data is the K2 algorithm [9]. Given a database D, this algorithm searches for the Bayesian network structure G^* with maximal $\Pr(G^*|D)$, where $\Pr(G|D)$ is the probability of network structure G given the database D. Since $\Pr(G_1|D)/\Pr(G_2|D) = \Pr(G_1,D)/\Pr(G_2,D)$ the authors look for a method to compute $\Pr(G,D)$. Let $V(G)$ be a set of n discrete variables, where a variable $V_i \in V(G)$ has r_i possible value assignments v_{ik} $k=1, \dots, r_i$. Let D be a database of m cases, where each case contains a value assignment for each variable in $V(G)$. Let G denote a directed acyclic graph representing the structure of a Bayesian network containing just the variables in $V(G)$, and let

GPr be the associated set of conditional probability distributions. Each node $V_i \in V(G)$ has a set of parents $\pi(V_i)$. Let w_{ij} denote the jth unique instantiation of $\pi(V_i)$ relative to D. Suppose there are q_i unique instantiations of $\pi(V_i)$. Define N_{ijk} to be the number of cases in D in which variable V_i has the value v_{ik} and $\pi(V_i)$ is instantiated as w_{ij} . Let

$$N_{ij} = \sum_{k=1}^{r_i} N_{ijk} \quad (1)$$

Given a Bayesian network structure G, assuming that the cases occur independently and the conditional probability density function $f(\text{GPr} | G)$ is uniform, then it follows that

$$\Pr(G, D) = \Pr(G) \prod_{i=1}^n \prod_{j=1}^{q_i} \frac{(r_i - 1)!}{(N_{ij} + r_i - 1)!} \prod_{k=1}^{r_i} N_{ijk}! \quad (2)$$

The K2 algorithm looks for a network structure G that maximizes $\Pr(G,D)$. In particular, assuming that an ordering on the variables is available and that all structures are equally similar, it adopts a greedy method for maximizing $\Pr(G,D)$. This method consists in, for every node V_i , searching for the set of parent nodes that maximizes the function:

$$g(i, \pi(V_i)) = \prod_{j=1}^{q_i} \frac{(r_i - 1)!}{(N_{ij} + r_i - 1)!} \prod_{k=1}^{r_i} N_{ijk}! \quad (3)$$

K2 starts by assuming that a node lacks parents, after which in every step it adds incrementally the parent whose addition mostly increases $g(i, \pi(V_i))$. K2 stops adding parents to a node when the addition of a single parent cannot increase $g(i, \pi(V_i))$.

A pseudo code representation of K2 algorithm is shown in Figure 1.

```

1. for i=1 to n
2. {
3.    $\pi(V_i)=0$ ;
4.   repeat
5.   {
6.     select  $V_j \in \{V_1, \dots, V_{i-1}\} - \pi(V_i)$  that maximizes
7.        $g(i, \pi(V_i) \cup \{V_j\})$ ;
8.      $\Delta = g(i, \pi(V_i) \cup \{V_j\}) - g(i, \pi(V_i))$ ;
9.     if  $\Delta > 0$  then  $\pi(V_i) = \pi(V_i) \cup \{V_j\}$ ;
10.  } until  $\Delta < 0$  or  $\pi(V_i) = \{V_1, \dots, V_{i-1}\}$ ;
11. }
```

Figure 1: Pseudo code representation of the K2 algorithm.

3 Association Rules

Association rules [1] describe correlation of events and can be considered as probabilistic rules. Events are “correlated” when they are frequently observed together. Good examples from real life are databases of sales transactions. In this case the aim is to find the items that are usually bought together, information that is used to develop successful marketing strategies. Consider a database D as the one given in Table 1.

Table 1: Example of a table of knowledge extraction

Attribute ₁	Attribute ₂	...	Attribute _n
Value _{1,1}	Value _{1,2}	...	Value _{1,n}
...
Value _{m,1}	Value _{m,2}	...	Value _{m,n}

An association rule is a rule of the form:

$$A_1=v_{A_1}, A_2=v_{A_2}, \dots, A_j=v_{A_j} \Rightarrow B_1=v_{B_1}, B_2=v_{B_2}, \dots, B_k=v_{B_k}$$

where $A_1, A_2, \dots, A_j, B_1, B_2, \dots, B_k$ are attribute names and $v_{A_1}, v_{A_2}, \dots, v_{A_j}, v_{B_1}, v_{B_2}, \dots, v_{B_k}$ are values such that v_{A_i} (v_{B_h}) belongs to the domain of the attribute A_i (B_h).

More formally, an association rule can be defined as follows.

An *item* is a literal of the form $Attribute_i=v_{Attribute_i}$ where $v_{Attribute_i}$ belongs to the domain $Attribute_i$. Let I be the set of all the possible items. A *transaction* T is a record of the database.

An *itemset* X is a set of items, i.e. it is a set X such that $X \subseteq I$. We say that a transaction T contains an itemset X if $X \subseteq T$ or, alternatively, if T satisfies all the literals in X .

The *support* of an itemset X (indicated by $support(X)$) is the fraction of transactions in D that contain X . The support of the opposite of an itemset (indicated by $support(!X)$) is the fraction of transactions in D that do not contain X . Therefore, $support(!X) = 1 - support(X)$.

An *association rule* is an implication of the form $X \Rightarrow Y$, where X and Y are itemsets and $X \cap Y \neq \emptyset$. For an association rule $X \Rightarrow Y$ we consider the following parameters:

- The *support* of $X \Rightarrow Y$ (indicated by $support(X \Rightarrow Y)$) is $support(X \cup Y)$.
- The *leverage* [26] of $X \Rightarrow Y$ (indicated by $leverage(X \Rightarrow Y)$) is given by $leverage(X \Rightarrow Y) = support(X \cup Y) - support(X) \times support(Y)$. (this parameter is similar to the Absolute Confidence Difference to Prior defined in [6]). It can assume positive and negative values.

Since $support(X)$ can be interpreted as $Pr(X)$, the leverage can be interpreted as $Pr(X, Y) - Pr(X) \times Pr(Y)$. Therefore the leverage is closer to 0 the more X and Y are statistically independent from each other.

In this paper we consider association rules where both X and Y contain a single item. In this way the leverage of the rule can be interpreted as a measure of the dependence between the two items contained respectively in X and Y .

4 K2-rules algorithm

In this section we describe the K2-rules algorithm which improves the K2 algorithm described in Section 3 by exploiting the leverage parameter of association rules. The K2 algorithm, in order to work, requires the total ordering of the nodes. K2 has a high computational cost and produces a significant number of extra arcs in the learned network.

The high computational cost is due to function (3) (see Section 2) which requires many computational resources especially for nodes characterized by a great number of parents.

The extra arc problem arises especially when the network is characterized by a lot of root nodes (nodes without parents). During network learning, the algorithm tries to add parents to each of these nodes until it maximizes function (3). The algorithm will add at least one arc to a root node because the value of the heuristics for this new structure is always better than the value of the previous structure.

The new proposed approach considers all the association rules containing an item in the body and an item in the head. In order to obtain the leverage of these rules, we do not employ an algorithm to learn association rules (such as APRIORI [2]), but we consider all the possible two items rules and for each we compute the

leverage. The K2-rules algorithm first computes, for each stochastic variable V_i , the maximum and the minimum of the leverage of the association rules that have an item that refers to V_i . Let $\text{MaxLev}(V_i)$ and $\text{MinLev}(V_i)$ be these figures.

Then K2-rules finds the minimum of all the $\text{MaxLev}(V_i)$ and the maximum of all the $\text{MinLev}(V_i)$. Let MaxLeverage and MinLeverage be these figures. Using these parameters, K2-rules deletes nodes from the list of possible parents of a node Q (those that precede it in the order). These parameters are used as thresholds for considering a couple of nodes statistically independent: if the leverage of a rule involving the two nodes is between MinLeverage and MaxLeverage , then the two variables are considered independent. Therefore the node that precedes the other in the given order can be removed from the list of parents of the other node.

```

Given the set of network nodes V and the set of
learned association rules AR:
1 For i=1 to n
2 {
3 Selects the subset AR(Vi) of association rules
  from AR which involve Vi;
4 Find the minimum and maximum value of
  leverage of rules in AR(Vi) and call them
  MinLev(Vi) and MaxLev(Vi);
5 }
6 Finds the global minimum and maximum for all
  the network nodes:
7 MinLeverage=maxVi ∈ V{MinLev(Vi)} and
  MaxLeverage=minVi ∈ V{MaxLev(Vi)}
8 For i=1 to n
9 {
10 π(Vi)=0;
11 Computes FindAllowableParents(Vi) or
  FindAllowableParents_All(Vi) which return a list
  AllowableParents of allowable nodes
12 Repeat
13 {
14 select Vj ∈ AllowableParents - πi that
  maximizes g(i, π(Vi) ∪ {Vj});
15 Δ = g(i, π(Vi) ∪ {Vj}) - g(i, π(Vi));
16 if Δ > 0 then π(Vi) = π(Vi) ∪ {Vj};
17 } until Δ < 0 or π(Vi) = {V1, ..., Vi-1};
18 }

```

Fig. 2. K2-rules algorithm

This method is implemented by the function `FindAllowableParents` that, given a node, returns the set of allowable parents. We have also considered a more conservative method to remove nodes from the set of allowable parents.

```

Given the ordered list of network nodes L, a node
and the set of learned association rules AR, the
FindAllowableParents(Node Q) function works as
follows:

```

- 1 Associates to Q a list `AllowableParents` of possible parents composed by the nodes that precede Q in the list L ;
- 2 Selects the subset $\text{AR}(Q)$ of association rules from AR which involve Q ;
- 3 For each node P in `AllowableParents`, if at least a rule R in $\text{AR}(Q)$ exists that involves P and Q , and that has $\text{MinLeverage} < \text{leverage}(R) < \text{MaxLeverage}$ then remove P from `AllowableParents`;
- 4 Return `AllowableParents`.

Fig. 3 `FindAllowableParents` function

This method is implemented by the function `FindAllowableParents_All`. This function deletes a node P from the list of parents of a node Q only if all the rules involving Q and P are such that their leverage is between MinLeverage and MaxLeverage . The K2-rules algorithm is described in pseudo code in Figure 2. The functions `FindAllowableParents` and

`FindAllowableParents_All` are presented in Figures 3 and 4 respectively.

```

Given the ordered list of network nodes L, a node
and the set of learned association rules AR, the
FindAllowableParents_All(Node Q) function works
as follows:

```

- 1 Associates to Q a list `AllowableParents` of possible parents composed by the nodes that precede Q in the list L ;
- 2 Selects the subset $\text{AR}(Q)$ of association rules from AR which involve Q ;
- 3 For each node P in `AllowableParents`, if all the rules $R \in \text{AR}(Q)$ which involve P and Q have $\text{MinLeverage} < \text{leverage}(R) < \text{MaxLeverage}$ then remove P from `AllowableParents`;
- 4 Return `AllowableParents`.

Fig. 4. `FindAllowableParents_All` function

5 Experimental comparisons

We compared K2 and K2-rules on the three different Bayesian networks:

- “Visit to Asia”: a network for a fictitious medical example about whether a patient has tuberculosis, lung cancer or bronchitis, depending on their X-ray, dyspnea, visit-to-Asia and smoking status. It has 8 nodes and 8 arcs, and is described in [16].

- “Car_diagnosis”: a network to diagnose the reason why a car does not start, based on spark plugs, headlights, main fuse, etc. It has 18 nodes and 20 arcs, and is described in [18];
- “ALARM”: ALARM stands for “A Logical Alarm Reduction Mechanism”. This is a medical diagnostic network to monitor patients. It is a nontrivial network with 8 diagnoses, 16 findings and 13 intermediate variables (37 nodes and 46 arcs), and is described in [4].
- “Boelarge92”: A subjective belief network for a particular scenario of neighborhood events, that shows how even distant concepts have some connection. It has 24 nodes and 35 arcs. It is described in [5];

This tool, given the structure and the CPTs of a Bayesian network is able to generate automatically a dataset of N examples. Each experiment was conducted by first generating a dataset from one of the networks above and then trying to learn back the network using K2, K2-rules using FindAllowableParents (K2-r-FAP), K2-rules using FindAllowableParents_All (K2-r-FAPA) and SLA. The learned networks are compared to the original network in Tables 2, 3, 4 and 5. For each algorithm we indicate: the numbers of missing and extra arcs (MA and EA, respectively); the Log Score (LS) indicating the number of computations of the function $g(i, \pi(V_i))$; the number of computation of N_{ijk} (NN). The last two parameters represent the computational resources needed by the K2, K2-r-FAP and K2-r-FAPA algorithms.

The dataset of examples used for rule learning has been obtained with the NETICA tool [18].

Table 2. Comparison between K2, K2-r-FAP and K2-r-FAPA for the Visit-to-Asia network.

Data Set	K2				K2-r-FAP				K2-r-FAPA				SLA	
	MA	EA	LS	NN	MA	EA	LS	NN	MA	EA	LS	NN	MA	EA
1000	0	4	66	2280	1	2	52	1980	1	2	52	1980	1	0
5000	0	2	61	1608	0	1	50	1416	0	1	50	1416	1	0
10000	0	1	57	1224	1	0	40	900	1	2	44	960	1	0
20000	0	1	57	1224	1	0	39	852	1	1	47	1020	1	0

Table 3. Comparison between K2, K2-r-FAP and K2-r-FAPA for the Car diagnosis network.

Data Set	K2				K2-r-FAP				K2-r-FAPA				SLA	
	MA	EA	LS	NN	MA	EA	LS	NN	MA	EA	LS	NN	MA	EA
1000	3	9	396	21705	4	6	270	16710	4	9	331	20679	8	0
5000	1	7	395	26817	2	2	214	17631	1	7	325	24375	6	0
10000	1	7	395	26817	2	1	201	16716	1	5	334	24837	7	0
20000	1	7	395	26817	1	0	206	17490	1	4	311	23745	6	0

Table 4. Comparison between K2, K2-r-FAP and K2-r-FAPA for the ALARM network.

Data Set	K2				K2-r-FAP				K2-r-FAPA				SLA	
	MA	EA	LS	NN	MA	EA	LS	NN	MA	EA	LS	NN	MA	EA
1000	3	13	1793	252120	23	3	143	11919	2	1	937	155178	37	37
5000	1	11	1771	204462	27	4	119	20163	1	0	685	91572	37	36
10000	1	11	1771	204462	22	3	148	23782	1	0	771	95235	37	35

Table 5. Comparison between K2, K2-r-FAP and K2-r-FAPA for the Boelarge network.

Data Set	K2				K2-r-FAP				K2-r-FAPA				SLA	
	MA	EA	LS	NN	MA	EA	LS	NN	MA	EA	LS	NN	MA	EA
1000	10	3	554	1194	10	1	196	4056	10	1	196	4056	13	1
5000	7	2	585	14244	7	0	172	4068	7	0	176	4140	12	0
10000	7	4	601	15732	8	2	199	5136	7	2	202	5284	13	0
20000	7	4	615	17172	8	3	161	4320	7	3	268	4500	12	0

Analyzing these experimental results we can observe that K2-r-FAP and K2-r-FAPA have a number of missing arcs comparable with that of K2 (apart from K2-r-FAP applied to ALARM) but have lower numbers of extra arcs. Moreover, the computational costs of both K2-rules algorithms are significantly lower than those required for K2. In particular, K2-r-FAP is more selective than K2-r-FAPA so it requires the least amount of resources.

Comparing K2-r-FAP and K2-r-FAPA with SLA we can observe that: SLA underestimates the probabilistic relations so it produces a high number of missing arcs while K2-r-FAP and K2-r-FAPA overestimates the probabilistic relations so they produce a low number of missing arcs but introduces some extra arcs. The total number of erroneous (missing and extra) arcs of SLA is higher or than the one of the new K2-rules algorithms (especially K2-r-FAPA) except for Visit to Asia and the first two datasets of Car_diagnosis.

6 Related Works

In this paper we present an approach for exploiting parameters associated to association rules in order to improve the performance of an algorithm for learning Bayesian networks. Such an approach is not limited to the K2 algorithm only. In fact, in [15], we have applied the same methodology to the SLA algorithm [7] that is based on an information theory approach rather than on a search and score methodology. In particular, we have exploited association rules parameters for improving the drafting phase of SLA. This phase is devoted to learn an initial sketch of the network structure.

In [15] we have used the parameters leverage, conviction, lift, Pearson X^2 and Cramer index. We have tested our algorithm (called BNL-rules) with each parameter on four networks together with SLA. In each test three

dimensions of the database were considered: 5.000, 20.000 and 100.000. Of all the algorithm tested, BNL-rules with Pearson X^2 gave the best results outperforming SLA in five cases in terms of number of missing and extra arcs, equating it in six cases and having a lower performance in only one case.

7 Conclusions

In this work we describe a method for improving K2, one of the most known algorithm for learning Bayesian network by exploiting association rules parameters.

The K2 algorithm starts by assuming that a node has no parents, after which in every step it incrementally adds the parent whose addition mostly increases the probability of the resulting structure. K2 stops adding parents to the nodes when the addition of a single parent cannot increase the probability of the resulting network structure given the data.

In this work, we propose a method to improve the K2 algorithm, reducing the set of allowable parents from which the algorithm selects actual parents and avoiding extra arc insertions. This new methodology uses data mining techniques, and in particular the computation of association rules parameters from a database of examples, in order to learn the structure of a Bayesian network. Association rules describe correlation of events, and are characterized by several parameters that can be used in structure learning. We have presented the K2-rules algorithm (K2 with association rules) that exploits the leverage parameter of association rules in order to improve the performance the K2 algorithm.

Experiments discussed in the paper have shown that the proposed approach solves the problem of extra arcs and also notably reduces the computational cost. They also showed the validity of the new algorithms with respect to SLA.

In future, we plan to compare K2-rules with MDL [25] and other Bayesian network learning algorithms.

Acknowledgments

We would like to thank Andrea Stambazzi and Rossella Bozzini for their help with the experiments. This work is partially funded by the Information Society Technologies programme of the European Commission under the IST-2001-32530 SOCS project.

References

- [1] R. Agrawal, T. Imielinski, A. Swami (1993). Mining association rules between sets of items in large databases. In: Buneman, P., Jajodia, S. (eds.): *Proceedings of the 1993 ACM SIGMOD International Conference on Management of Data*, ACM Press, volume 22(2), pages 207-216
- [2] R. Agrawal, R. Srikant (1994). Fast Algorithms for Mining Association Rules. In: Bocca, J.B., Jarke, M., Zaniolo, C. (eds.): *Proceedings of the 20th International Conference on Very Large Data Bases (VLDB'94)*. Morgan Kaufmann, volume ISBN 1-55860-153-8, pages 487-499
- [3] Akaike, H. (1974). A new Look at Statistical Model Identification. *IEEE Trans. Automatic Control*, volume 19, pages 716-723
- [4] I.A. Beinlich, H.J. Suermondt, R.M. Chavez, G.F. Cooper (1989). The ALARM monitoring system: A case study with two probabilistic inference techniques for belief networks. In: Hunter, J. (ed.): *Proceedings of the Second European Conference on Artificial Intelligence in Medicine (AIME 89)*. Springer, Berlin, pages 247-256
- [5] Boerlage, Brent (1992). Link Strength in Bayesian Networks. *MSc Thesis*, Dept. Computer Science, Univ. of British Columbia, BC.
- [6] C. Borgelt. <http://fuzzy.cs.uni-magdeburg.de/~borgelt/doc/apriori/apriori.html#diff>
- [7] J. Cheng, R. Greiner, J. Kelly, D. Bell, W. Liu (2002). Learning Bayesian networks from data: An information-theory based approach, *Artificial Intelligence*, volume 137, pages 43-90
- [8] C.K. Chow, C.N. Liu (1968). Approximating discrete probability distributions with dependence trees. *IEEE Transactions on Information Theory*, volume 14, pages 462-467
- [9] G. Cooper, E. Herskovits (1992). A Bayesian method for the induction of probabilistic networks from data. *Machine Learning*, volume 9, pages 309-347
- [10] P. Giudici (2001). Data mining - Metodi statistici per le applicazioni aziendali. *McGraw-Hill*
- [11] D. Heckerman, D. Geiger, D. Chickering (1995). Learning Bayesian Networks: the combination of knowledge and statistical data. *Machine Learning*, volume 20, pages 197-243
- [12] D. Heckerman (1999). Tutorial on learning in Bayesian networks. In: Jordan, M. (ed.): *Learning in Graphical Models*. MIT Press
- [13] E.H. Herskovits (1991). Computer-based probabilistic-network construction. Doctoral Dissertation, Medical Informatics, Stanford University
- [14] J.H. Kim, J. Pearl (1983). A Computational Model for Combined Causal and Diagnostic Reasoning in Inference Systems. In *Proceedings of the Eighth International Joint Conference on Artificial Intelligence (IJCAI83)*. pages 190-193
- [15] E. Lamma, F. Riguzzi, A. Stambazzi, S. Storari (2003). Improving the SLA algorithm using association rules. *Proceedings of the Eighth Italian Congress on Artificial Intelligence (AI*IA2003)*, LNAI 2829, Springer Verlag
- [16] S.L. Lauritzen, D.J. Spiegelhalter (1988). Local computations with probabilities on graphical structures and their application to expert systems. *J. Royal Statistics Society B*, volume 50(2), pages 157-194
- [17] D. Madigan, A. Raftery (1994). Model Selection and Accounting for Model Uncertainty in Graphical Models Using

- Occam's Window. *J. Am. Statist. Association*, volume 89, pages 1535-1546
- [18] Netica, <http://www.norsys.com>
- [19] J. Pearl (1998). Probabilistic reasoning in intelligent systems: networks of plausible inference. *Morgan Kaufmann*
- [20] J. Rissanen (1987). Stochastic Complexity (with discussion). *J. Roy. Statist. Soc. B*, volume 49, pages 223-239
- [21] G. Schwarz (1978). Estimating the Dimension of a Model. *Annals of Statistics*, volume 6, pages 461-464
- [22] M. Sigh, M. Valtorta (1993). Construction of Bayesian Network Structures from Data: a Brief Survey and an Efficient Algorithm. In: Heckerman, D., Mamdani, A. (eds.): *Proceedings of the Conference on Uncertainty in Artificial Intelligence*. Morgan Kaufmann, pages 259-265
- [23] P. Spirtes, C. Glymour, R. Scheines (1993). Causation, Prediction, and Search. *Lecture Notes in Statistics*. Springer
- [24] P. Spirtes, C. Glymour, R. Scheines (1991). An algorithm for fast recovery of sparse causal graphs. *Social Science Computer Review*, volume 9, pages 62-72
- [25] J. Suzuki (1999). Learning Bayesian Belief Networks Based on the MDL principle: An Efficient Algorithm Using the Branch and Bound Technique. *IEICE Transactions on Communications Electronics Information and Systems*.
- [26] Weka Manual Page http://www.cs.waikato.ac.nz/~ml/weka/doc_gui/weka.asociations.ItemSet.html