

Initialization Heuristics for Greedy Bayesian Network Structure Learning

Walter Perez Urcia

and

Denis Deratani Mauá

Instituto de Matemática e Estatística, Universidade de São Paulo, Brazil
wperez@ime.usp.br, denis.maua@usp.br

Abstract. A popular and effective method for learning Bayesian network structures is to perform a greedy search on the space of variable orderings followed by an exhaustive search over the restricted space of compatible parent sets. Usually, the greedy search is initialized with a randomly sampled order. In this article we develop heuristics for producing informed initial solutions to order-based search motivated by the Feedback Arc Set Problem.

Categories and Subject Descriptors: I.2.6 [Artificial Intelligence]: Learning

Keywords: Bayesian networks, machine learning, local search

1. INTRODUCTION

Bayesian Networks are space-efficient representations of multivariate probability distributions over the attributes of a data set. They are defined by two components: (i) a directed acyclic graph (DAG), where the nodes are the attributes of the data set and the edges encode the (in)dependence relationships among the attributes; and (ii) a collection of local conditional probability distributions of each attribute given its parents. More formally, a Bayesian network specification contains a DAG $G = (V, E)$, where $V = \{X_1, X_2, \dots, X_n\}$ is the set of attributes, and a collection of conditional probability distributions $P(X_i | Pa_G(X_i))$, $i = 1, \dots, n$, where $Pa_G(X_i)$ is the set of attributes that are parents of X_i in G . This definition shows that the number of numerical parameters (i.e., local conditional probability values) grows exponentially with the number of parents (in-degree) of a node. A Bayesian network induces a joint probability distribution such that

$$P(X_1, X_2, \dots, X_n) = \prod_{i=1}^n P(X_i | Pa_G(X_i)).$$

Hence, Bayesian networks with sparse DAGs are succinctly represent joint probability distributions.

A common approach to learning Bayesian networks from a given data set is the search-and-score approach, which associates every graph structure with a polynomial-time computable score value; we then search for structures with high score values. The score of a structure is usually a balance between the probability of observing the data set and the complexity of the model (i.e., the number of parameters). Some examples are the Bayesian Information Criterion (BIC) [Cover and Thomas 1991], the Minimum Description Length (MDL) [Lam and Bacchus 1994] and the Bayesian Dirichlet score (BD) [Heckerman et al. 1995].

Score-based Bayesian network structure learning from data is a NP-hard problem [Chickering et al. 2004], even when the in-degree (i.e., maximum number of parents) of the graph is bounded. For this reason, the commonest approach to solve the problem is to use local search methods such as searching over the space of structures, searching over Markov equivalence classes, and searching over the space of topological orders. Local search methods are well-known to be vulnerable to poor local maxima unless a strategy for avoiding low score regions is used. One such strategy is the use of non-greedy heuristics that allow moving for lower value solutions during search in order to escape local maxima. Another strategy is to use *good initialization heuristics* that attempt to start the search in regions of high local maxima. Most often, a simple uniformed random generation of initial solutions is adopted.

In this article we propose new heuristics for generating good initial solutions to be fed into local search methods for Bayesian network structure learning. We focus on order-based methods, but our technique can be exploited by any local search procedure for learning Bayesian network structures. Our heuristics are motivated by solutions of the Feedback Arc Set Problem (FASP), which is the problem of transforming a cyclic directed graph into a DAG. Our experiments show that using these new methods improves the quality of order-based local search.

The article is structured as follows: We begin in Section 2 explaining the Greedy Search algorithm. In Section 3, we describe the new algorithm for generating initial solutions. Section 4 shows the experiments using both approaches and comparing them (in time and scoring) with multiple data sets. Finally, in Section 5 we give some conclusions about the new methods.

2. LEARNING BAYESIAN NETWORKS

In this section, we define mathematically the learning Bayesian network problem and then reduce it to one formula in order to be used as a heuristic score. Finally we explain some popular methods to solve them.

2.1 Definition of the problem

The problem of learning the structure of a Bayesian network is stated as: Given a training data set D , select the network (DAG) G that maximizes the scoring function $sc(G, D)$. A general definition of a scoring function is as follows:

$$sc(G, D) = F(G, D) - \varphi(N) \times P(G),$$

where N is the size of the data set D , $F(G, D)$ is a data fitness function (i.e., how well the model represents the observed data), $\varphi(N)$ is a non-decreasing function of data size and $P(G)$ measures the model complexity of G . For instance, the Bayesian information criterion (BIC) is defined as

$$BIC(G, D) = LL(G) - \frac{\log N}{2} size(G),$$

where

$$LL(G) = \sum_{i=1}^n \sum_k \sum_j N_{ijk} \log \frac{N_{ijk}}{N_{ij}},$$

is the data loglikelihood

$$size(G) = \sum_{i=1}^n (|\Omega_i| - 1) \prod_{X_j \in Pa(X_i)} |\Omega_j|,$$

represents the “size” of a model with structure G , N is the number of instances in the data set, n is the number of attributes, N_{ijk} the number of instances where attribute X_i takes values k , and its parents take the j th configuration (in some arbitrary ordering), and similarly for N_{ij} , and Ω_i is the set of possible values for the attribute X_i .

Most scoring functions, BIC included, are *decomposable*, meaning that they can be written as a sum of local scoring functions $sc(X_i, Pa(X_i))$. Under the assumption of *score decomposability*, the Bayesian network structure learning problem is defined as finding

$$G = \max_{\{Pa(X_i)\}_i, G \text{ is acyclic}} \sum_{i=1}^n sc(X_i, Pa(X_i)).$$

In words, the structure learning problem is to select for each attribute the set parent attributes that maximize the local score while ensuring that the graph is acyclic. Exhaustively searching for the optimal parent set $Pa(X_i)$ for each attribute takes $O(2^n)$ time. To avoid the exponential blow-up we constrain the maximum number of parents to a value d , hence reducing the search space to 2^d ; this means $|Pa(X_i)| \leq d$ for all attributes. Typical values for d are 3, 4 and 5, depending on the cardinality of variables. Additional speed-up can be achieved by using pruning schemes that remove suboptimal subsets of parent sets without resorting to inspection.

2.2 Greedy Search Algorithm

The Greedy Search algorithm is a popular heuristic method used to find an approximation for a solution of learning Bayesian network problem using the concept of neighborhood between solutions. Depending on the focus of the algorithm this could be classified as Equivalence-based, Structure-based, Order-based, etc. Algorithm 1 shows a general pseudocode for this algorithm.

Algorithm 1: Greedy Search

```

1   $L = Initial\_Solution(X_1, \dots, X_n, conf)$ 
2   $best\_score = score(L)$ 
3  For a number of iterations  $K$  or until  $L$  converges do
4       $best\_neighbor = find\_best\_neighbor(L)$ 
5      if  $score(best\_neighbor) > score(L)$  then
6           $L = best\_neighbor$ 
7  Return  $L$ 
```

The main idea of this algorithm is to generate an initial solution based on the *conf* parameter (for instance, *conf* can say that the solution has to be random generated). After that, for a number of iterations K or until the solution converges (the best solution's score does not improve) explore the search space and selects the best neighbor of the best solution until that moment. Then calculate its score in order to know if it is a better solution. Finally, return the best solution when the stop condition holds.

As mentioned before, there are lot of different approaches using this algorithm, but all of them only change lines 1, 4 and 2 depending on its own way to generate an initial solution, its search space and the scoring function used, respectively.

For example, in the structure-based version the initial solution is a network, possibly random generated, and its neighbors are other networks that only differs by one arc. This means that an arc can be inserted, deleted or inverted its direction to obtain a new neighbor.

In the same way, the equivalence-based version [Chickering and Meek 2004] uses the equivalence relation between networks to obtain new neighbors. This relation two networks L_1 and L_2 are equivalent if they have the same probability distribution.

Finally, a Order-based Greedy Search is a popular and effective solution to the problem of learning the structure of a Bayesian network. Algorithm 2 shows its pseudocode, where the function *swap* (in line 6) swaps the values $L[i]$ and $L[i + 1]$ in the order L .

Algorithm 2: Order-based Greedy Search

```

1   $L = Get\_Order(X_1, \dots, X_n, conf)$ 
```

```

2  best_score = score(L)
3  For a number of iterations K or until L converges do
4      current_sol = L
5      For each i = 1 to n - 1 do
6          Li = swap(L, i, i + 1)
7          if score(Li) > score(current_sol)
8              current_sol = Li
9          if score(current_sol) > score(L) then
10             L = current_sol
11  Return L

```

So we need to calculate the score for an order of the attributes and we know from subsection 2.1 that, in general, to obtain the score for a network is an NP-hard problem because of the exponential number of possible sets of parents (even in the bounded case). But if we know an order of the attributes, the problem can be reduced as follows:

$$G = \max_G \sum_{i=1}^n sc(X_i, Pa(X_i)) = \sum_{i=1}^n \max_{P \subseteq \{X_j < X_i\} : |P| \leq d} sc(X_i, P),$$

where $X_j < X_i$ means that X_j appears before X_i in the order of the attributes. This means that we can maximize the score for every attribute independently the other ones.

In section 4, the Order-based Greedy Search algorithm will be used for learning structure of Bayesian networks using multiple data sets.

3. GENERATING INFORMED INITIAL SOLUTIONS

As stated in Section 1, the solutions proposed in this article are motivated by the Feedback Arc Set Problem (FASP) that will be explained at the start of this section. After that, we explain every new method to generate an initial solution for Order-based Greedy Algorithm.

3.1 Feedback Arc Set Problem

The generic problem is stated as: Given a graph $G = (V, E)$, a set $F \subseteq E$ is called the Feedback Arc Set (FAS) if every cycle of G has at least one edge in F . In other words, F is the edge set that if you put it off the graph G , it becomes a directed acyclic graph (DAG).

There are some variants of this problem, but If we focus on unweighted undirected graphs and we want to find a FAS with minimum size, the problem becomes NP-hard, but there are some approximation algorithms like the one that is below.

- (1) Select a random node $X \in V$
- (2) Do a depth-first search (DFS) on G using X as root
- (3) Put all the back edges encountered by the DFS in F . A back edge is the one that makes a cycle with size greater than two
- (4) Return F

In the same way, another NP-hard problem is to find a minimum weight FAS from a weighted undirected graph G , but again we have some approximation algorithms in order to solve this problem.

- (1) Negate all edge weights
- (2) Find the minimum spanning tree (MST) M performing Kruskal algorithm
- (3) Edges that are not in M are included in F

(4) Return F

The variants of the problem where the graph is directed (with weighted or unweighted edges) will not be explained in this article.

3.2 Informed initial solutions

From section 2.2, we already know the generic Greedy Search algorithm and the more popular approaches using it. It can be noticed that a very important step in the algorithm is in the first line, where an initial solution is generated because if we have a good one, the solution will converge faster. So if we now focus only in the Order-based one showed in Algorithm 2, we first explain the common approach and then we will propose two new methods for generating initial solutions using information from data motivated by the solutions in 3.1.

3.2.1 Random solution. This is the most common and easy-to-implement approach used. An easy way to shuffle a list L with size n is the Fisher-Yates algorithm, showed in Algorithm 3.

Algorithm 3: Fisher-Yates shuffle algorithm

```

1      For  $i$  in  $n - 1$  to 1 do
2           $j = \text{random}(0, i)$ 
3          Swap  $L[i]$  and  $L[j]$ 
4      Return  $L$ 
```

3.2.2 Using FAS with unweighted edges. Remembering the definition of a Bayesian network given in 2.1 we have:

$$G = \max_G \sum_{i=1}^n sc(X_i, Pa(X_i))$$

Although the Bayesian network G has the maximum sum of the scores of every field X_i given their parents ($Pa(X_i)$), this does not imply that every field has their best parents because the graph built with them will probably not be a network as it can contain cycles. Besides that, calculating the best parents is computationally costly unless we bound the size of the parent set as mentioned in 2.1. Knowing those concepts and the approximation algorithm for unweighted FAS problem, we can propose the following algorithm:

- (1) For each node X_i : Find the $Pa(X_i)$ set with maximum score
- (2) Build the graph S using all the relationships $X_j \implies X_i, X_j \in Pa(X_i)$
- (3) Choose a node X_k from S
- (4) Start with an empty graph G
- (5) Perform a DFS on using X_k as root with the following characteristics:
 - Add the direct edge e that encounters to G , unless e builds a cycle
 - If e builds a cycle, then add e to the feedback arc set F
- (6) For each edge in F : Invert its direction
- (7) Add edges in F to the DAG G
- (8) Return the topological order of G

This solution is similar to the one for the unweighted FAS problem using the graph S , except that we use the tree built by the DFS to get a network and then add the Feedback Arc Set with their directions inverted. Finally, as we want a field order, we return the topological order of the network G .

3.2.3 *Using FAS with weighted edges.* In this version, the edges have its weight defined as follows:

$$W_{ji} = sc(X_i, P) - sc(X_i, \{X_j\}),$$

where W_{ji} is the weight of the edge that goes from X_j to X_i and P is the best parents set for X_i . This formula represents the cost of getting X_j out of the set P and it is always a positive number. When the value is so small, it means the parent X_j could be more important. On the contrary if the weight is so big, this means that X_j is not so important in P .

The algorithm for getting an initial solution is very similar to the previous one and it uses the main idea of the weighted FAS approximation algorithm.

- (1) For each node X_i : Find the $Pa(X_i)$ set with maximum score
- (2) Build the graph S using all the relationships $X_j \implies X_i, X_j \in P = Pa(X_i)$ and put its score as $-W_{ji}$
- (3) Make the graph S undirected and perform Kruskal algorithm for getting the MST M
- (4) Choose a node X_k from M
- (5) Start with an empty graph G
- (6) Perform a DFS on M and add the directed edges on G
- (7) Return the topological order of G

In step 3, we have to make the graph S undirected because Kruskal algorithm can be used only in undirected weighted graphs. Finally, we perform a DFS to get a directed tree with root X_k and we only have to return its topological order.

These new methods for generating initial solutions will be used in next section to learn Bayesian networks with multiple data sets.

4. EXPERIMENTS AND RESULTS

After implementing all the approaches in section 3.2, some experiments were done in order to compare them. First, we show the setup for all experiments and then show the final results.

4.1 Experiments Setup

Table I shows the characteristics of each data set used in the experiments. Also, the values for each

Dataset	n (#attributes)	N (#instances)
Alarm	37	1000
Census	15	45222
Epigenetics	30	72228
Image	20	2310
Letter	17	20000
Mushroom	23	8124
Sensors	25	5456
SteelPlates	28	1941

Table I: Data sets characteristics

general parameter are given below.

- Maximum number of parents (d): 3
- Number of iterations in Greedy Search (K): 100
- Scoring function used: Bayesian Criterion Information (BIC)

Finally, as in subsections 3.2.2 and 3.2.3 we have to choose a field as a root for performing the DFS, we have n possible different DAGs and then n possible different orders and use each one as an initial solution for Greedy Search. Also, for the random version, we generate n different random orders.

4.2 Results

After running the three versions explained in 3.2 for each data set, we compare the score for the best network found, the CPU Time to generate the initial solution and the CPU Time for getting the final network. Tables II and III show them respectively.

Dataset	Random Sol.	Unweighted Sol.	Weighted Sol.
Alarm	37	1000	0
Census	15	45222	0
Epigenetics	30	72228	0
Image	20	2310	0
Letter	17	20000	0
Mushroom	23	8124	0
Sensors	25	5456	0
SteelPlates	28	1941	0

Table II: Maximum cost using each approach

Dataset	Random Sol.	Unweighted Sol.	Weighted Sol.
Alarm	37	1000	0
Census	15	45222	0
Epigenetics	30	72228	0
Image	20	2310	0
Letter	17	20000	0
Mushroom	23	8124	0
Sensors	25	5456	0
SteelPlates	28	1941	0

Table III: CPU Time (in seconds) for initializing and total using each approach

Finally, Figure 1 shows the converge curve for each data set using each approach.

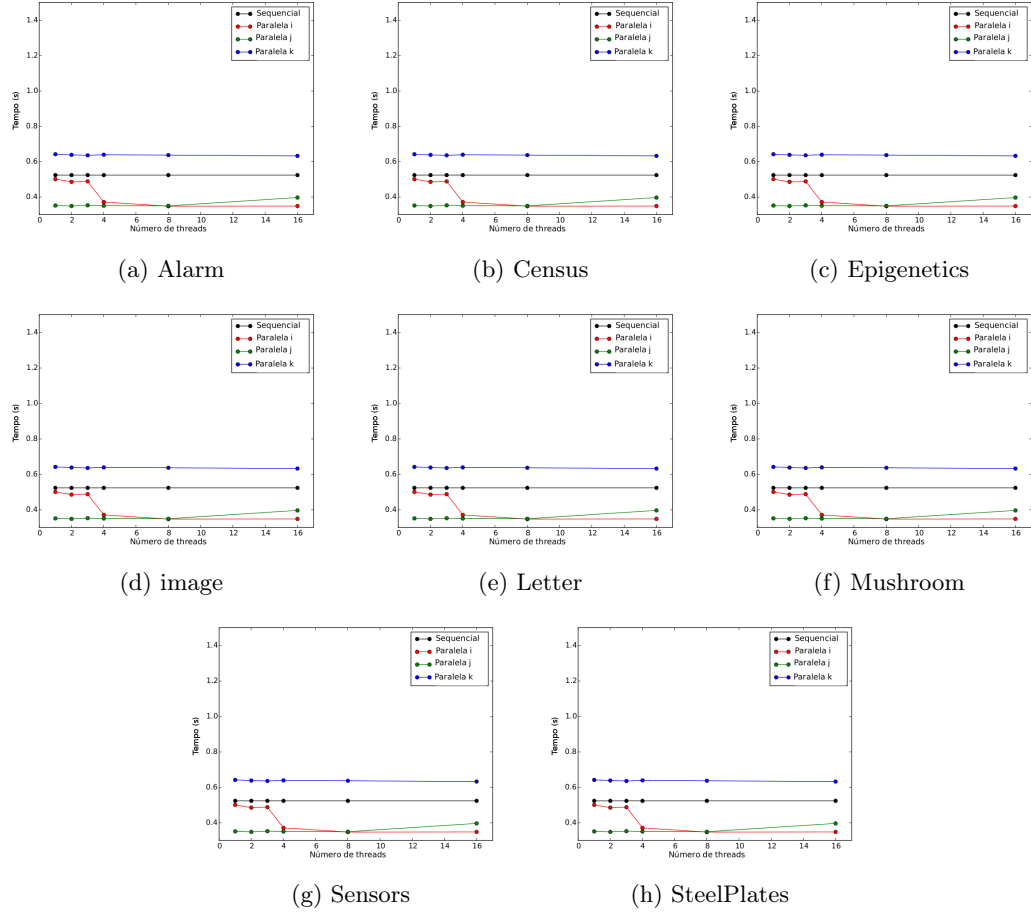


Fig. 1: Converge curve for datasets

5. CONCLUSIONS AND FUTURE WORK

We conclude that:

- Using the random initial solution take more iterations to find a converged solution than the new methods and also more CPU time
- The unweighted approach gives best networks than the other ones
- Initial solutions using the unweighted approach are the most time consuming to be generated and the random ones the least.
- A higher number of attributes implies more CPU time for the new methods with a significant difference compared to the random approach
- The new methods explained in this article could be used in other types of Greedy Search doing some modifications
- The scoring function could be changed to another one in order to compare robustness between them and BIC score
- Experiments could be done with bigger datasets (more than 20 fields) in order to compare the efficiency of the methods

REFERENCES

- CHICKERING, D. M., HECKERMAN, D., AND MEEK, C. Large-Sample Learning of Bayesian Networks is NP-Hard. *Journal of Machine Learning Research* 5 (1): 1287–1330, 2004.
- CHICKERING, D. M. AND MEEK, C. Finding Optimal Bayesian Networks. *Journal of Machine Learning Research*, 2004.
- COVER, T. M. AND THOMAS, J. A. *Elements of Information Theory*. Wiley-Interscience, 1991.
- HECKERMAN, D., GEIGER, D., AND CHICKERING, D. Learning Bayesian Networks: The Combination of Knowledge and Statistical Data. *Journal of Machine Learning Research* 20 (MSR-TR-94-09): 197–243, 1995.
- LAM, W. AND BACCHUS, F. Learning Bayesian Belief Networks. An approach based on the MDL principle. *Computational Intelligence* 10 (4): 31, 1994.