

# 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 approach 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 on data sets without missing values.

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 complex multivariate probability distributions [Jensen 2001]. They are defined by two components: (i) a directed acyclic graph (DAG) encoding the (in)dependence relationships among the variables in the model; and (ii) a collection of local conditional probability distributions of each variable given its parents.

Manually specifying a Bayesian network is a difficult task, and practitioners often resort to “learning” the model from data. A common approach to learning a Bayesian network consists of associating every DAG with a polynomial-time computable score value and searching for structures with high score values [Cooper and Dietterich 1992; Lam and Bacchus 1994; Margaritis 2003; Tessyer and Koller 2005]. The score value of a structure usually rewards structures that assign high probability of observing the data set (i.e., the data likelihood) and penalizes 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]. An alternative approach is to learn the DAG by multiple conditional independence hypothesis testing [Spirtes and Meek 1995; Cheng et al. 2002]. Although both approaches can recover the true DAG (if one exists) given infinite data and computational resources, testing for independence introduces a lot of false positives and it is often followed by a score-based approach [Tsamardinos et al. 2006].

Score-based Bayesian network 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 most common approach is to resort local search methods that find an approximate solution [H. Friedman and Peér 1999; Chickering 2002]. A popular and very effective method for learning

Bayesian networks is to perform a local search on the space of topological orderings [Tessier and Koller 2005]. The search is usually initialized with an ordering sampled uniformly at random from the space of orderings. This can make the search converge to a poor local optima unless more sophisticated techniques are employed [Elidan et al. 2002], which can add significant computational overhead. An alternative solution is to initialize the search in high-scoring regions.

In this work we design two new heuristics for generating good initial solutions to order-based Bayesian network structure learning. The first heuristic follows the observation that only orderings consistent with a relaxed version of the problem (in which cycles are permitted) can lead to an optimal structure. Although this heuristic biases the search away from regions which are *guaranteed* to be sub-optimal, it generates orderings with equal probability in any other region. Our second heuristic refines the first one by selecting high scoring orderings among the ones that are consistent with the relaxed version solution. We do this by reducing the problem to a variant 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 rest of this paper is structured as follows: we begin in Section 2 explaining greedy search approaches to learning Bayesian networks. Then in Section 3 we describe the new algorithms for generating initial solutions. Section 4 shows the experiments using both approaches and comparing them (in scoring and number of iterations needed) with multiple data sets. Finally, in Section 5 we give some conclusions about the new methods.

## 2. LEARNING BAYESIAN NETWORKS

In this section, we formally define the score-based approach learning of Bayesian networks, and review some of the most popular techniques for solving the problem.

### 2.1 Definition of the problem

A Bayesian network specification contains a DAG  $G = (V, E)$ , where  $V = \{X_1, X_2, \dots, X_n\}$  is the set of (discrete) variables, 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 variables 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 (assuming the values are organized in tables). A Bayesian network induces a joint probability distribution over all the variables through the equation  $P(X_1, X_2, \dots, X_n) = \prod_{i=1}^n P(X_i | Pa_G(X_i))$ . Hence, Bayesian networks with sparse DAGs succinctly represent joint probability distributions over many variables.

A *scoring function*  $sc(G)$  assigns a real-value to any DAG indicating its goodness in representing a given data set.<sup>1</sup> Most scoring functions can be written in the form  $sc(G) = F(G) - \varphi(N) \times P(G)$ , where  $N$  is the number of records in the data set  $D$ ,  $F(G)$  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 example, the Bayesian information criterion (BIC) is defined as  $BIC(G) = LL(G) - \frac{\log N}{2} size(G)$ , where  $LL(G) = \sum_{i=1}^n \sum_k \sum_j N_{ijk} \log \frac{N_{ijk}}{N_{ij}}$  is the data loglikelihood,  $size(G) = \sum_{i=1}^n (|\Omega_i| - 1) \prod_{X_j \in Pa(X_i)} |\Omega_j|$  is the “size” of a model with structure  $G$ ,  $n$  is the number of attributes on  $D$ ,  $N_{ijk}$  the number of instances where attribute  $X_i$  takes its  $k$ th value, and its parents take the  $j$ th configuration (for some arbitrary fixed ordering of the configurations of the parents’ values), and similarly for  $N_{ij}$ , and  $\Omega_i$  is the set of possible values for the attribute  $X_i$ . Most commonly used scoring functions, BIC included, are *decomposable*, meaning that they can be written as a sum of local scoring functions:  $sc(G) = \sum_i sc(X_i, Pa(X_i))$ . Another property often

<sup>1</sup>The dependence of the scoring function on the data set is usually left implicitly, as for most of this explanation we can assume a fixed data set. We assume here that the dataset contains no missing values.

satisfied by scoring functions is *likelihood equivalence*, which asserts that two structures with same loglikelihood also have the same score [Chickering and Meek 2004]. Likelihood equivalence is justified as a desirable property, since two structures that assign the same loglikelihood to data cannot be distinguished by the data alone. The BIC scoring function satisfies likelihood equivalence.

Given scoring function  $sc(G)$ , the score-based Bayesian network structure learning problem is to compute the DAG

$$G^* = \arg \max_{G: G \text{ is a DAG}} sc(G). \quad (1)$$

Provided the scoring function is decomposable, we can obtain an upper bound on the value of  $sc(G^*)$  by computing  $sc(\bar{G})$ , where

$$\bar{G} = \arg \sum_i \max_{Pa(X_i)} sc(X_i, Pa(X_i)) \quad (2)$$

is the directed graph where the parents  $Pa(X_i)$  of each node  $X_i$  are selected so as to maximize the local score  $sc(X_i, Pa(X_i))$ . We call the parents of a variable in  $\bar{G}$  the *best parent set* (for  $X_i$ ). Note that  $\bar{G}$  usually contains cycles, and it is thus not a solution to equation 1.

## 2.2 Greedy Search Approaches

Greedy Search is a popular approach used to finding an approximate solution to equation (1). The method relies on the definition of a neighborhood space among solutions, and on local moves that search for an improving solution in the neighborhood of an incumbent solution. Different neighborhoods and local moves give rise to different methods such as Equivalence-based, Structure-based, and Order-based methods. Algorithm 1 shows a general pseudocode for this approach.

Algorithm 1: Greedy Search

```

1  GreedySearch( Dataset  $D$  ) : return a BN  $G$ 
2     $G = \text{Initial\_Solution}(X_1, \dots, X_n)$ 
3    For a number of iterations  $K$ 
4       $best\_neighbor = \text{find\_best\_neighbor}(G)$ 
5      if  $\text{score}(best\_neighbor) > \text{score}(G)$  then
6         $G = best\_neighbor$ 
7    Return  $G$ 
```

The main idea of the approach is to start with an initial solution (e.g., a randomly generated one), and for a number of iterations  $K$ , explore the search space by selecting the best neighbor of the incumbent solution. Additionally, an early stop condition can be added to verify whether the algorithm has reached a local optimum (i.e., if no local move can improve the lower bound). Several methods can be obtained by varying the implementation of lines 2, 4 and 5, which specify how to generate an initial solution, what the search space is and what the scoring function is, respectively.

**2.2.1 Structure-based.** One of earliest approaches to learning Bayesian networks was to perform a greedy search over the space of DAGs, with local moves being the operations of adding, removing or reverting an edge, followed by the verification of acyclicity in the case of edge addition [Cooper and Dietterich 1992; Grzegorzczuk and Husmeier 2008]. The initial solution is usually obtained by randomly generating a DAG, using one of the many methods available in the literature [Ide and Cozman 2002; Melançon and Philippe 2004].

**2.2.2 Equivalence-based.** An alternative approach is to search within the class of score-equivalent DAGs. This can be efficiently achieved when the scoring function is likelihood equivalent by using pDAGs, which are graphs that contain both undirected and directed edges (but no directed cycles) with the property that all orientations of a pDAG have the same score. In this case, greedy search

operates on the space of pDAGs, and the neighborhood is defined by addition, removal and reversal of edges, just as in structure-based search [Chickering 1996; 2002].

**2.2.3 Order-based.** Order-Based Greedy Search is a popular and effective approach, which is based on the observation that the problem of learning a Bayesian network can be written as

$$G^* = \arg \max_{<} \max_{G \text{ consistent with } <} \sum_{i=1}^n sc(X_i, Pa(X_i)) = \arg \max_{<} \sum_{i=1}^n \max_{P \subseteq \{X_j < X_i\}} sc(X_i, P), \quad (3)$$

which means that if an optimal ordering over the variables is known, an optimal DAG can be found by maximizing the local scores independently [Heckerman et al. 1995; H. Friedman and Peér 1999; Tessier and Koller 2005]. This can be made efficiently if we assume  $G^*$  is sparse, which is true for many scoring functions [de Campos and Ji 2011].

Order-Based Search starts with a topological ordering  $L$ , and greedily moves to an improving ordering by swapping two adjacent attributes in  $L$  if any exists. Algorithm 2 shows a pseudocode for the method. The function *swap* in line 6 swaps the values  $L[i]$  and  $L[i + 1]$  in the order  $L$  to obtain a neighbor of the incumbent solution.

Algorithm 2: Order-Based Greedy Search

```

1  OrderBasedGreedySearch( Dataset  $D$  ) : return a BN
2     $L = \text{Get\_Order}(X_1, \dots, X_n)$ 
3    For a number of iterations  $K$ 
4       $current\_sol = L$ 
5      For each  $i = 1$  to  $n - 1$  do
6         $L_i = \text{swap}(L, i, i + 1)$ 
7        if  $score(L_i) > score(current\_sol)$ 
8           $current\_sol = L_i$ 
9        if  $score(current\_sol) > score(L)$  then
10          $L = current\_sol$ 
11    Return  $network(L)$ 
```

The standard approach to generate initial solutions is to sample a permutation of the attributes uniformly at random by some efficient procedure such as the Fisher-Yates algorithm [Knuth 1998]. While this guarantees a good coverage of the search space when many restarts are performed, it can lead to poor local optima. In the next section, we propose new strategies to informed generation of topological orderings to be used as initial solutions in Order-Based search.

### 3. GENERATING INFORMED INITIAL SOLUTIONS

As with most local search approaches, the selection of a good initial solution is crucial for avoiding convergence to poor local maxima in Order-Based Learning. Traditionally, this is attempted by randomly generating initial solutions (i.e., a node ordering) in order to cover as much as possible of the search space. In this section, we devise methods that take advantage of the structure of the problem to produce better initial solutions.

#### 3.1 DFS-based approach

We can exploit the information provided by the graph  $\overline{G}$  (defined in equation 2) to reduce the space of topological orderings and avoid generating orderings which are guaranteed sub-optimal. Assume the best parent sets are unique, and consider a pair of nodes  $X_i, X_j$  in  $\overline{G}$  such that  $X_j$  is parent of  $X_i$  but there is not arc from  $X_i$  into  $X_j$ . Then, no optimal ordering can have  $X_i$  preceding  $X_j$  (this can easily be shown by contradiction). Hence, only topological orderings consistent with  $\overline{G}$  are potential candidates for optimality, and this number can be much smaller than the full space of orderings. To

see this clearly, consider Figure 1 which shows a possible graph  $\overline{G}$  and the corresponding consistent orderings. As can be noticed we have 14 consistent orderings out of  $4! = 24$  possible topological orders. This difference is likely to increase as the number of variables increases.

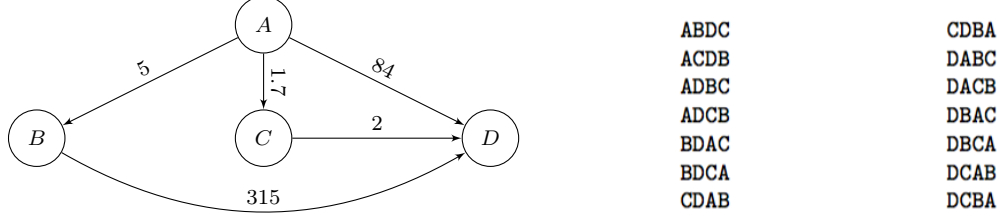


Fig. 1: A an example of a graph  $\overline{G}$  and its consistent topological orderings

Taking into consideration the previous analysis, we propose the following algorithm to generate initial solutions. Take as input the graph  $\overline{G}$  and mark all nodes as unvisited. While there is an unvisited node, select an unvisited node  $X_i$  uniformly at random and add to the list the nodes visited by a depth-first search (DFS) tree rooted at  $X_i$ . Finally, return  $L$ , an ordering of the nodes.

### 3.2 FAS-based approach

The DFS approach can be seen as removing edges from  $\overline{G}$  such as to make it a DAG (more specifically, a tree), and then extracting a consistent topological ordering. That approach hence considers that all edges are equally relevant in terms of avoiding poor local maxima. We can estimate the arguably relevance of an edge  $X_j \rightarrow X_i$  by

$$W_{ji} = sc(X_i, Pa^*(X_i)) - sc(X_i, Pa^*(X_i) \setminus \{X_j\}), \quad (4)$$

where  $Pa^*(X_i)$  denotes the best parent set for  $X_i$  (i.e., its parents in  $\overline{G}$ ). The weight  $W_{ji}$  represents the cost of removing  $X_j$  from the set  $Pa^*(X_i)$  and it is always a positive number because  $Pa(X_i)$  maximizes the score for  $X_i$ . A small value means that the parent  $X_j$  is not very relevant to  $X_i$  (in that sense), while a large value denotes the opposite. For instance, in the weighted graph  $\overline{G}$  in Figure 1, the edge  $C \rightarrow D$  is less relevant than the edges  $A \rightarrow D$ , which in turn is less relevant than the edge  $B \rightarrow D$ .

The main idea of our second heuristic is to penalize orderings which violate an edge  $X_i \rightarrow X_j$  in  $\overline{G}$  by their associated cost  $W_{ij}$ . We then wish to find a topological ordering of  $\overline{G}$  that violates the least cost of edges. Given a directed graph  $G = (V, E)$ , a set  $F \subseteq E$  is called a Feedback Arc Set (FAS) if every (directed) cycle of  $G$  contains at least one edge in  $F$ . In other words,  $F$  is an edge set that if removed makes the graph  $G$  acyclic [Demetrescu and Finocchi 2003]. If we assume that the cost of an ordering of  $\overline{G}$  is the sum of the weights of the violated (or removed) edges, we can formulate the problem of finding a minimum cost ordering of  $\overline{G}$  as a Minimum Cost Feedback Arc Set Problem (min-cost FAS): given the weighted directed graph  $\overline{G}$  with weights  $W_{ij}$  given by equation (4), find a FAS  $F$  such that

$$F = \min_{G-F \text{ is a DAG}} \sum_{X_i \rightarrow X_j \in E} W_{ij}. \quad (5)$$

Even though the problem is NP-hard, there are efficient and effective approximation algorithms like the one described in Algorithm 3 [Demetrescu and Finocchi 2003].

#### Algorithm 3: FAS approximation

```

1  MinimumCostFAS( Graph  $G$  ) : Return FAS  $F$ 
2   $F = \text{empty set}$ 
```

```

3      While there is a cycle  $C$  on  $G$  do
4           $W_{min}$  = lowest weight of all edges in  $C$ 
5          For each edge  $(u, v) \in C$  do
6               $W_{uv} = W_{uv} - W_{min}$ 
7              If  $W_{uv} = 0$  add to  $F$ 
8          For each edge in  $F$ , add it to  $G$  if does not build a cycle
9      Return  $F$ 

```

We can now describe our second heuristic for generating initial solutions, based on the minimum cost FAS problem: take the weighted graph  $\overline{G}$  with weights  $W_{ij}$  as input, and find a min-cost FAS  $F$ ; remove the edges in  $F$  from  $\overline{G}$  and return a topological order of the obtained graph  $\overline{G} - F$  (this can be done by performing a DFS starting with root nodes).

#### 4. EXPERIMENTS, RESULTS AND DISCUSSION

In order to evaluate the quality of our approaches, we learned Bayesian networks using Order-based greedy search and different initialization strategies from several data sets commonly used for benchmarking. The names and relevant characteristics of the data sets<sup>2</sup> used are shown in Table I, where the density of a graph is defined as the ratio of the number of edges and the number of nodes. For

Dataset	n (#attributes)	N (#instances)	Density of $\overline{G}$
Census	15	30168	2.85
Letter	17	20000	2.41
Image	20	2310	2.45
Mushroom	23	8124	2.91
Sensors	25	5456	3.00
SteelPlates	28	1941	2.18
Epigenetics	30	72228	1.87
Alarm	37	1000	1.98
Spectf	45	267	1.76
LungCancer	57	27	1.44

Table I: Data sets characteristics

each dataset we performed 1000 runs of Order-Based Greedy Search with a limit of 3 parents ( $d = 3$ ) and 100 iterations ( $K = 100$ ), except for the LungCancer dataset where only 100 runs were performed due to the limited computational resources. We used the BIC score and found the best parent sets for a given ordering by exhaustive search.

We compared our proposed initialization strategies, which we call DFS- and FAS-based, against the standard approach of randomly generating an order (called Random). For each strategy, we compared the best score obtained over all runs (Best score), the average initial score (i.e., the score of the best DAG consistent with the initial ordering), the average best score (i.e., the average of the scores of the local searches) and the average number of iterations that local search took to converge. The results are shown in Table II. The results show that in most of the datasets with less than 25 attributes, the Random strategy finds the highest-scoring networks over all runs, even though it finds worse networks on average. The best initial solutions are found by the FAS-based strategy followed by the DFS-based strategy. For datasets with more than 25 variables, Random is less effective in finding high-scoring networks, except for the LungCancer (which has very little data). These results suggest that more informed approaches to generating initial orderings might be more effective in high dimensionality domains, or when the number of restarts is limited e.g. for computational reasons. The proposed

<sup>2</sup>These datasets were extracted from <http://urlearning.org/datasets.html>

Dataset	Approach	Best Score	Avg. Initial Score	Avg. Best Score	Avg. It.
Census	Random	<b>-212186.79</b>	-213074.18 $\pm$ 558.43	-212342.26 $\pm$ 174.21	7.26 $\pm$ 2.90
	DFS-based	-212190.05	-212736.80 $\pm$ 379.96	-212339.83 $\pm$ 152.26	5.90 $\pm$ 2.61
	FAS-based	-212191.64	<b>-212287.99 <math>\pm</math> 92.54</b>	<b>-212222.12 <math>\pm</math> 70.99</b>	<b>3.28 <math>\pm</math> 1.67</b>
Letter	Random	-138652.66	-139774.54 $\pm$ 413.74	-139107.13 $\pm$ 329.15	6.07 $\pm$ 2.50
	DFS-based	-138652.66	-139521.38 $\pm$ 396.61	<b>-138999.84 <math>\pm</math> 310.06</b>	5.75 $\pm$ 2.35
	FAS-based	-138652.66	<b>-139050.43 <math>\pm</math> 70.55</b>	-139039.26 $\pm$ 87.97	<b>2.24 <math>\pm</math> 0.96</b>
Image	Random	<b>-12826.08</b>	-13017.13 $\pm$ 44.35	-12924.24 $\pm$ 41.39	7.59 $\pm$ 2.71
	DFS-based	-12829.10	-12999.09 $\pm$ 38.56	-12921.13 $\pm$ 37.88	7.10 $\pm$ 2.47
	FAS-based	-12829.10	<b>-12930.63 <math>\pm</math> 20.83</b>	<b>-12882.30 <math>\pm</math> 26.43</b>	<b>5.05 <math>\pm</math> 1.72</b>
Mushroom	Random	<b>-55513.38</b>	-58450.72 $\pm$ 1016.54	-56563.84 $\pm$ 616.59	7.59 $\pm$ 2.76
	DFS-based	<b>-55513.38</b>	-58367.11 $\pm$ 871.25	-56472.72 $\pm$ 546.19	7.75 $\pm$ 2.58
	FAS-based	-55574.71	<b>-56450.49 <math>\pm</math> 154.54</b>	<b>-56198.66 <math>\pm</math> 174.64</b>	<b>4.65 <math>\pm</math> 1.63</b>
Sensors	Random	<b>-62062.13</b>	-63476.33 $\pm$ 265.46	-62726.60 $\pm$ 251.26	9.22 $\pm$ 2.94
	DFS-based	-62083.21	-63392.60 $\pm$ 255.90	-62711.50 $\pm$ 257.79	9.65 $\pm$ 3.12
	FAS-based	-62074.88	<b>-62530.26 <math>\pm</math> 133.44</b>	<b>-62330.94 <math>\pm</math> 121.82</b>	<b>5.17 <math>\pm</math> 2.24</b>
SteelPlates	Random	-13336.14	-13566.50 $\pm$ 65.80	-13429.13 $\pm$ 52.14	8.96 $\pm$ 3.43
	DFS-based	<b>-13332.91</b>	-13572.77 $\pm$ 81.12	-13432.30 $\pm$ 57.57	9.30 $\pm$ 3.38
	FAS-based	-13341.73	<b>-13485.26 <math>\pm</math> 38.27</b>	<b>-13397.08 <math>\pm</math> 29.53</b>	<b>7.77 <math>\pm</math> 2.24</b>
Epigenetics	Random	-56873.76	-57722.30 $\pm$ 228.44	-57357.60 $\pm$ 222.12	5.89 $\pm$ 2.67
	DFS-based	<b>-56868.87</b>	<b>-57615.36 <math>\pm</math> 189.17</b>	<b>-57308.93 <math>\pm</math> 165.18</b>	6.42 $\pm$ 2.47
	FAS-based	<b>-56868.87</b>	-57660.09 $\pm$ 146.45	-57379.59 $\pm$ 148.42	<b>5.33 <math>\pm</math> 2.28</b>
Alarm	Random	-13218.22	-13324.52 $\pm$ 30.49	-13245.43 $\pm$ 15.63	10.92 $\pm$ 3.24
	DFS-based	<b>-13217.97</b>	-13250.72 $\pm$ 17.70	-13236.71 $\pm$ 12.02	<b>4.32 <math>\pm</math> 2.32</b>
	FAS-based	-13220.55	<b>-13249.77 <math>\pm</math> 2.57</b>	<b>-13233.98 <math>\pm</math> 6.19</b>	6.34 $\pm$ 1.74
Spectf	Random	-8176.81	-8202.03 $\pm$ 5.23	-8189.69 $\pm$ 4.65	7.20 $\pm$ 2.17
	DFS-based	<b>-8172.37</b>	-8200.04 $\pm$ 4.08	-8187.29 $\pm$ 4.91	7.86 $\pm$ 2.49
	FAS-based	-8172.51	<b>-8176.98 <math>\pm</math> 2.01</b>	<b>-8176.07 <math>\pm</math> 2.05</b>	<b>2.27 <math>\pm</math> 1.11</b>
LungCancer	Random	<b>-711.23</b>	-723.79 $\pm$ 2.69	-718.03 $\pm$ 2.84	5.46 $\pm$ 1.78
	DFS-based	-711.36	-720.47 $\pm$ 2.51	<b>-715.29 <math>\pm</math> 1.86</b>	5.02 $\pm$ 1.50
	FAS-based	-711.39	<b>-716.13 <math>\pm</math> 0.89</b>	-715.67 $\pm$ 1.19	<b>2.73 <math>\pm</math> 1.79</b>

Table II: Best score obtained, Average initial score generated, Average best score obtained, Average number of iterations (Avg. It.) using each approach (best values in bold)

strategies are also more robust, which can be seen by the smaller variance of the average initial and best scores.

The results also suggest that the proposed strategies are more effective than Random in datasets for which the graph  $\bar{G}$  is sparser (smaller density), showing that pruning the space of orderings can be effective in those cases. The initial orderings provided by the proposed strategies speed up convergence of the local search, as can be seen by the smaller number of average iterations for those strategies in the table.

Overall, the new heuristics are able to improve the accuracy of Order-Based Greedy Search with only a small overhead. Although the differences observed in our experiments were small, we expect greater differences in domains of higher dimensionality.

## 5. CONCLUSIONS AND FUTURE WORK

Learning Bayesian networks from data is a notably difficult problem, and practitioners often resort to approximate solutions such as greedy search. The quality of the solutions produced by greedy approaches strongly depends on the initial solution. In this work, we proposed two new heuristics for producing topological orderings to be fed into Order-Based Greedy Bayesian network Structure Search methods. One is based on a Depth-First Search traversal of the (cyclic) graph obtained by greedily selecting the best parents for each variable; the other is based on finding an acyclic subgraph

of that same graph by solving a related minimum cost Feedback-Arc Set problem. Experiments with real-world datasets containing from 15 to 57 variables demonstrate that compared to the commonly used strategy of generating initial ordering uniformly at random the proposed heuristics lead to better solutions on average, and increase the convergence of the search with only a small overhead. Although the gains observed in our experiments are small, we expect larger differences for datasets with more variables. A follow-up work should verify this hypothesis.

Our proposed techniques could be adapted to generate initial solutions also for Structure- and Equivalence-based local search methods by returning directed acyclic graphs instead of node orderings. Another extension of this work is to employ the proposed heuristics in branch-and-bound solvers such as [de Campos and Ji 2011] for finding optimal solutions. These ideas are left as future work.

## REFERENCES

- CHENG, J., GREINER, R., KELLY, J., BELL, D., AND LIU, W. Learning bayesian networks from data: An information-theory based approach. *Artificial Intelligence* vol. 137, pp. 43–90, 2002.
- CHICKERING, D. M. Learning equivalence classes of Bayesian-network structures. *Conference on Uncertainty in Artificial Intelligence*, 1996.
- CHICKERING, D. M. Learning equivalence classes of Bayesian-network structures. *Journal of Machine Learning Research*, 2002.
- 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.
- COOPER, G. F. AND DIETTERICH, T. A Bayesian method for the induction of probabilistic networks from data. *Machine Learning*, 1992.
- COVER, T. M. AND THOMAS, J. A. *Elements of Information Theory*. Wiley-Interscience, 1991.
- DE CAMPOS, C. P. AND JI, Q. Efficient structure learning of Bayesian networks using constraints. *Journal of Machine Learning Research* vol. 12, pp. 663–689, 2011.
- DEMETRESCU, C. AND FINOCCHI, I. Combinatorial algorithms for feedback problems in directed graphs. *Information Processing Letters*, 2003.
- ELIDAN, G., NINIO, M., AND SCHUURMANS, N. F. D. Data perturbation for escaping local maxima in learning. *Proceedings of the National Conference on Artificial Intelligence*, 2002.
- GRZEGORCZYK, M. AND HUSMEIER, D. Improving the structure MCMC sampler for Bayesian networks by introducing a new edge reversal move. *Machine Learning*, 2008.
- H. FRIEDMAN, I. N. AND PEÉR, D. Learning Bayesian network structure from massive datasets: The “sparse candidate” algorithm. *Conference on Uncertainty in Artificial Intelligence* (15), 1999.
- 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.
- IDE, J. S. AND COZMAN, F. G. Random generation of Bayesian networks. vol. 2507, pp. 366–376, 2002.
- JENSEN, F. V. *Bayesian Networks and Decision Graphs*. Springer Science and Business Media, 2001.
- KNUTH. *The Art of Computer Programming 2*. Boston: Adison-Wesley, 1998.
- LAM, W. AND BACCHUS, F. Learning Bayesian belief networks. an approach based on the MDL principle. *Computational Intelligence* 10 (4): 31, 1994.
- MARGARITIS, D. Learning Bayesian network model structure from data, 2003.
- MELANÇON, G. AND PHILIPPE, F. Generating connected acyclic digraphs uniformly at random. *Information Processing Letters* 90 (4): 209–213, May, 2004.
- SPIRITES, P. AND MEEK, C. Learning Bayesian networks with discrete variables from data. *Proceedings of the 1st International Conference on Knowledge Discovery and Data Mining*, 1995.
- TESSYER, M. AND KOLLER, D. Ordering-based search: A simple and effective algorithm for learning Bayesian networks. *Proceedings of the Conference in Uncertainty in Artificial Intelligence*, 2005.
- TSAMARDINOS, I., BROWN, L. E., AND ALIFERIS, C. F. The max-min hill-climbing Bayesian network structure learning algorithm. *Machine Learning* vol. 65, pp. 31–78, 2006.