

**NOTICE WARNING CONCERNING COPYRIGHT RESTRICTIONS:**

The copyright law of the United States (title 17, U.S. Code) governs the making of photocopies or other reproductions of copyrighted material. Any copying of this document without permission of its author may be prohibited by law.

# Tractable Structural Learning of Large Bayesian Networks from Sparse Data

*Anna Goldenberg and Andrew W. Moore*  
anyaOcs.emu.edu, awmOcs.emu.edu

April 28, 2004

CMU-CALD-04-1033

School of Computer Science  
Carnegie Mellon University  
Pittsburgh, PA 15213

## Abstract

This paper addresses three questions. Is it useful to attempt to learn a Bayesian network structure with hundreds of thousands of nodes? How should such structure search proceed practically? The third question arises out of our approach to the second: how can Frequent Sets (Agrawal et al., 1993), which are extremely popular in the area of descriptive data mining, be turned into a probabilistic model? Large sparse datasets with hundreds of thousands of records and attributes appear in social networks, warehousing, supermarket transactions and web logs. The complexity of structural search made learning of factored probabilistic models on such datasets unfeasible. We propose to use Frequent Sets to significantly speed up the structural search. Unlike previous approaches, we not only cache  $n$ -way sufficient statistics, but also exploit their local structure. We also present an empirical evaluation of our algorithm applied to several massive datasets.

University Libraries  
Carnegie Mellon University  
Pittsburgh, PA 15213-3890

**Keywords:** Bayesian networks/graphical models, statistical learning, Bayes Net structure learning

# 1 Introduction

Bayesian Networks have been successfully applied in many areas such as pharmaceutical, decision making by doctors, air control, marketing, etc. Structural learning of Bayesian Networks is usually a desirable but costly operation. In some domains it is possible to collect expert knowledge to manually create a structure for a Bayes Net. However, social networks, warehousing data, or supermarket purchasing records may contain hundreds of thousands of attributes. Providing expert Bayes Net structure in such cases is cumbersome if not impossible, even if as in the case with many of those domains the events are choices of very small subsets of the large pool of available entities. The complexity of existing algorithms for structural search prevents Bayes Net learning on datasets of that size.

This paper provides an algorithm for tractable structural learning in Bayes Nets by exploring structures on the local level. We exploit the computational efficiency of Frequent Sets for gathering statistics that are most likely to be useful for structure search given the assumption of sparse data. We then give an efficient search algorithm to exploit these statistics for creating the global Bayes Net.

Usage of Bayesian Networks to represent expression of genes based on the activity of their regulators (in practice approximated by protein activity levels) is well motivated by Friedman (2004). He suggests that the structure of the network is of its own importance, since it may provide information about gene's regulators.

Another field that has received increasing attention in the last few years is recommender systems. A lot of online systems such as Amazon provide suggestions of what might appeal to the user based on user's other preferences. The use of Bayesian models in this domain has been demonstrated by (Breese et al., 1998). Often the goal of recommender systems is to predict which are the most likely items that the user would buy next. An example of answering analogous query using Bayes Nets built by our algorithm is presented in Section 6.2.

The idea of representing social networks as people connected by directed arrows has been explored in social science domain for almost 70 years (Moreno & Jennings, 1938). Initially analyzed networks were on the order of 10s of nodes. However, improvements in data collection and especially the birth of online communities made it necessary to look at much larger networks. For example, livejournal (an online blog community) contains over 50,000 users (Shklovsky, personal communication). Usage of graphical models in this domain has become increasingly popular, due to their robustness to noise.

Note that even though in every domain described above, there are potentially tens or hundreds of thousands variables, each variable interacts only with a chosen few, resulting in very sparse datasets. For example, it has been long known that connectivity in social networks follows power law distribution (Barabasi, 2002), i.e. there are very few people who are connected to many people, most of them are only interacting with a very small group relative to the number of people available in the net.

In fact, studies in the gene expression data and social networks in particular suggest that correlations of entities on the local level are very important and in fact they are what makes up the global network (Friedman, 2004; Breiger, 2003). So, along with being computationally practical Bayesian Networks created by our algorithm have a very natural motivation

stemming from those important domains.

Also, we claim that for the descriptive data mining community, this algorithm may help to answer an old question: “what should you do with Frequent Sets once you’ve found them?”.

We provide results on sparse massive datasets showing practical training times, and in many cases superior ability to model the joint distribution in comparison with direct extensions of traditional structure search algorithms on large data. We also qualitatively and empirically show that sparse data particularly with social net characteristics are modelled better by going beyond information derived from pairwise co-occurrences.

## 2 Frequent Sets

Assume our training data is a collection of  $N$  records with  $M$  binary categorical attributes per record. Write  $x_{ij}$  as the value of the  $j$ th attribute of the  $i$ th record where  $1 \leq i \leq N$  and  $1 \leq j \leq M$ . We assume sparse data in which the vast majority of values in any dataset row or column are zero. Note we assume no missing values: all  $x_{ij}$  are observed.

Let the  $M$  attributes be represented by integers  $\{1, 2, \dots, M\}$ . Let the *co-occurrence frequency* of a set of attributes  $S \subseteq \{1, 2, \dots, M\}$  be the number of records in which all the attributes in  $S$  are simultaneously set to 1.

$$freq(S) = |\{i : \forall j \in S, x_{ij} = 1\}| \quad (1)$$

Given  $s \geq 1$  we say  $S$  is a *Frequent Set of  $m$  attributes* if  $S$  contains exactly  $m$  attributes and  $freq(S) \geq s$ . Threshold  $s$  is called *support* in the data mining literature. Given sparse data and a support  $s$  greater than about 3, it is surprisingly easy to compute all Frequent Sets (Agrawal & Srikant, 1994). There is an abundance of literature on Frequent Sets as their collection is an essential part of the association rules algorithms (Agrawal et al., 1993; Agrawal & Srikant, 1994; Han & Kamber, 2000) widely used in commercial data mining.

There are multiple references to Frequent Sets in the area of modelling sparse datasets as well (Mannila & Toivonen, 1996; Chickering & Heckerman, 1999; Pavlov et al., 2003; Hollmen et al., 2003). This is not surprising, since sparseness implies very few co-occurrences between items. In fact, most items do not co-occur with each other, hence we expect the majority of the counts in the pairwise marginals to be 0. Therefore, it is natural to assume that the Frequent Sets contain most of the essential information about the whole dataset.

(Chickering & Heckerman, 1999) propose and show how to use an efficient sparse representation for several classes of machine learning algorithms including structure initialization for Bayes Nets. We will therefore not focus on representational aspects of Frequent Sets.

This paper exploits previous research on the utilization of Frequent Sets for modelling of sparse datasets but takes a new perspective. Assuming that Frequent Sets comprise essential information about our data we propose to exploit them to find Bayes Net structures on the local level. To our knowledge, structures contained *within* Frequent Sets have not been previously used in order to improve the global model of data.

### 3 Algorithm Description

The simplest idea for exploiting Frequent Set information is to use frequent pairs. The only edges which we would consider including in the Bayes Net are those for which the source and destination attributes co-occur more than some support  $s$ . There are thus far fewer edges to consider than the full  $M(M - 1)$  possibilities ( $M$  is the number of attributes).

There are three problems with this idea.

- **Problem 1.** This method will not find edges that have negative correlations. For example, if attribute  $A$  is never positive when attribute  $B$  is positive then  $(A,B)$  will not be a frequent pair and so will not be considered.

**Solution.** Problem 1 is mitigated in two ways. First, under the assumption of sparse data there must necessarily be less evidence for a strong negative correlation as is shown in the Appendix. In fact, the structure scoring metric (e.g. BDeu) will be much higher for positively correlated entities. Secondly, attributes with high marginal positive values (where a negative correlation might be significant), will be accounted for at a later stage described in Section 4.

- **Problem 2.** Some items that do co-occur might be independent. This is especially likely with promiscuous attributes that occur frequently by themselves and thus could co-occur just by chance.

**Solution.** The solution to problem 2 is to screen all frequent pairs before allowing links between them into the pool of edges considered for the network. Only significantly correlated pairs become candidate edges. This greatly reduces the number of candidate edges.

- **Problem 3.** Restricting the search to frequent pairs can miss significant higher-order interactions. The appendix gives one example, but it is easy to imagine many cases in which co-occurrence of  $A$  and  $B$  is predictive of the occurrence of  $X$  and yet one or both of the  $A \rightarrow X$  and  $B \rightarrow X$  dependencies are not statistically significant.

**Solution.** This is solved by using higher-order Frequent Sets, as described in the following paragraphs.

#### 3.1 Screening the Frequent Sets.

We call the set of edges that will eventually be considered for addition into the Bayesian Network the *Edgedump*. Suppose we have a collection of Frequent Sets  $\{X : |X| = m, m \geq 2\}$ . First, we screen the pairs to find positive pairwise correlations. We add an edge between two variables to the *Edgedump* if and only if a significant correlation was found between the 2 variables in the pair. We then in turn screen for dependencies in Frequent Sets of size 3, 4, etc.

When does a Frequent Set  $X$  of size  $m > 2$  provide new information valuable for building a Bayes Net? It is possible that the dependencies of  $X$  are already well-explained by interactions of order less than  $m$ . For example, suppose attributes  $A$ ,  $B$  and  $C$  co-occur frequently,

but their co-occurrence is well explained by the local Bayesian Network DAG structure of  $A \leftarrow B \rightarrow C$ . In that case the two-way interactions will already explain all dependencies of  $X$ . In this case,  $X$  should not be added to the edgedump. In fact, only DAGs that contain a node with  $m - 1$  parents could be missed by considering only lower order interactions, as is shown in the Appendix for the case of triplets.

We implement a *Screening* test by searching over all possible DAG structures for  $X$  and finding whether the best BDeu-scoring structure has an  $m - 1$ -parent node (we call it an *m-way* interaction). We thus allow  $X$  to pass the screening test *if and only if*  $X$  is best explained by a local DAG structure containing an *m-way* interaction. If  $X$  passes the Screening test, all edges of the highest scoring DAG are added to the Edgedump.

Once the Edgedump is created, we prioritize the edges according to their strength, measured by the number of the *m-way* interactions in which they participate. We then create an empty (edgeless) global Bayesian network and iterate through the Edgedump contents, allowing each edge in turn to be added if and only if it improves BDeu and avoids cycles.

Table 1 contains the full description of the algorithm.

## 4 Addition of High Mutual Information Links

In the previous section we pointed out that Frequent Sets bias the algorithm in favor of interactions that cause co-occurrence (and thus positive correlation). Appendix 1 shows why, in the case of sparse data, positive correlations must be stronger than negative correlations, so in general we are not omitting the strongest correlations. There is, however, still a danger that if a few attributes are promiscuous (relatively high univariate marginal probability, though still very sparse), they could cause significant negative correlations that we could miss. Fortunately, such negative pairwise correlations can be detected cheaply using a technique from (Meila, 1999).

Let  $I_{AB}$  be the mutual information between two attributes. Meila showed that the mutual information can be calculated in a very efficient manner, particularly when dealing with discrete binary data. In fact, if the two variables have not co-occurred in the dataset, the formula simplifies even further:  $I_{AB}$  is directly proportional to the magnitude of  $A$  ( $N_A$ ) and  $B$  ( $N_B$ ) as shown in Equation (2). The full derivation is available in (Meila, 1999).

$$\begin{aligned}
 I_{AB} &= H_A + H_B - H_{AB} \\
 &= \frac{1}{N} [ -(N - N_A) \log(N - N_B) - (N - N_B) \\
 &\quad x \log(N - N_B) + (N - N_A - N_B) \\
 &\quad x \log(N - N_A - N_B) + N \log N ]
 \end{aligned} \tag{2}$$

Hence, to add high mutual information edges, we have to check entities that occur with high frequency. We reduce the total number of entities significantly by only considering ones that occurred more than  $s$  times in the dataset. This step is statistically justified because fewer occurrences mean lower possible mutual information. Table 2 describes the algorithm that augments a given Bayes Net with high mutual information (MI) edges.

Table 1: Screen-based Bayes Net Structure search (SBNS) algorithm

<b>algorithmSBNS</b>	
<b>input</b>	<b>K</b> - max Frequent Set size s - support
<b>output</b>	<i>BN</i> - Bayes Net
Also:	
<i>Ed</i>	Edgedump - a collection of directed edges represented as (source,dest,count)
<i>DS</i>	DAG storage
<ol style="list-style-type: none"> <li>1. <b>for</b> <math>k = 2 .. K</math></li> <li>2.   obtain counts for all Frequent Sets of size <math>k</math></li> <li>3.   <b>foreach</b> Frequent Set</li> <li>4.     find   best scoring DAG</li> <li>5.     <b>if</b> DAG contains a node that has <math>k - 1</math> parents</li> <li>6.       store DAG in <i>DS</i></li> <li>7.     <b>end foreach</b></li> <li>8. <b>end for</b></li> <li>9. <b>foreach</b> DAG in <i>DS</i></li> <li>10.   store all edges <math>\{source, dest, count++\}</math> in <i>Ed</i></li> <li>11.   order <i>Ed</i> in decreasing order of edge counts</li> <li>12. <b>foreach</b> edge <math>e \in Ed</math></li> <li>13.   <b>if</b> <math>e</math> doesn't form a cycle in <i>BN</i></li> <li>14.     <b>and</b> <math>e</math> improves <i>BDeu</i></li> <li>15.     add <math>e</math> to <i>BN</i></li> <li>16. <b>end foreach</b></li> <li>17. return <i>BN</i></li> </ol>	

We do not search the space of all edges to find edges with the highest mutual information. First of all, we sort entities in descending order of frequency. For each entity  $A^i = 1 \dots N_{>s}$ , where  $N_{>s}$  is the number of entities with support  $> s$ , we only consider  $\{B^j = i + 1 \dots N_{>s}\}$ , i.e. those entities that have occurred less frequently than  $A^i$ . If an edge  $e^{ij}$  has been rejected, then we move along the  $A$  list. This step is justified, because entities are sorted in descending order of frequencies, hence the mutual information between  $A^i$  and  $B^{j+i}$  is lower than between  $A^i$  and  $Z^j$ . Thus, the edge  $e_{A^i B^{j+i}}$  is even less likely to be added than  $e_{A^i B^j}$  - Empirical evidence shows that on average only 10% of  $N$  pairs are considered.



Table 2: Algorithm that augments Bayes Net  $BN$  with high MI edges

<b>algorithm</b>	AugmentWithMutualEdges
<b>input</b>	$BN$ - a Bayes Net $L$ - list of attributes with frequencies
<ol style="list-style-type: none"> <li>1. Sort <math>L</math> in decreasing order of frequencies</li> <li>2. <b>for</b> <math>u = 1 \dots  L  - 1</math></li> <li>3.     <math>v = u + 1</math>; <i>added.lag</i> = <i>TRUE</i></li> <li>4.     <b>while</b> <math>v &lt;  L </math> &amp; <i>added-flag</i></li> <li>5.         <b>if</b> net with <math>e_{uv}</math> score <math>&gt;</math> old net score</li> <li>6.             add <math>e_{uv}</math> to <math>BN</math></li> <li>7.         <b>else</b> try to add <math>e^{\wedge}</math> to <math>BN</math></li> <li>8.         <math>v = v + 1</math></li> <li>9.         <b>if</b> (edge not added) <i>added.flag</i> = <i>FALSE</i></li> <li>10.     <b>end while</b></li> <li>11. <b>end for</b></li> <li>12. return <math>BN</math></li> </ol>	

## 5 Additional possible postprocessing

### 5.1 Second Degree Separation Links

It is cheap to do an extra pass of edge-additions in which we iterate over all nodes in the network produced by the previous steps and attempt adding edges directly from the current node to its grandchildren.

### 5.2 Hillclimbing

One of the standard techniques to improve the score is hillclimbing as described in (Cooper & Herskovits, 1991). This technique improves the score by adding/removing/reversing arcs in a Bayes Net. The set of operations and edge selection procedure may differ between algorithms. Usually hillclimbing is performed in a *greedy* search way: at each step the existing model undergoes a modification/addition of a single edge. In order to pick the best edge we must look at  $O(N^2)$  possibilities. Since the number of nodes  $N$  prohibits us to perform even a linear search at each step, we use *random* hillclimbing in which at each step we choose edges randomly. Specifically, we roll a "3-sided" die with probabilities .8 for addition and .1 for deletion and arc-reversal, and then pick an edge at random to see whether performing the chosen operation improves the global score.

## 6 Evaluation

The evaluation uses BDeu score described in (Heckerman et al., 1995) and also presented here in equation 3 to compare results between different configurations of our algorithm and to the randomized hillclimbing as described in Section 5.2.

$$S(G, D) = \prod_{i=1}^n \prod_{j \in \text{pa}(i)} \sum_{r \in \{0,1\}} P(x_i = r | \text{states of } X_j) \quad \diamond$$

where  $i$  is the  $i$ th variable,  $x_i$  - states of the  $i$ th parent of  $x_i$ ,  $r$  - true/false (in our case of binary variables) states of  $X_i$ .

The datasets are listed in Table 3. Holdout testsets were used to evaluate overfitting as discussed in Section 6.2.3.

### 6.1 Datasets

The algorithm has been tested on several real life datasets<sup>1</sup> (sizes shown in Table 3).

1. The *Institute Data* is a set of records of collaborations between professors and students collected from publicly available web pages listed on Carnegie Mellon University Robotic Institute's web site.

---

<sup>1</sup>We will make the non-proprietary datasets available on the web upon publication

Table 3: Data sets and their sizes

Datasets	Entities	Records
Institute	456	1488
Drinks	136	4744
IMDB	100717	49298
Citeseer	104801	180395

2. The *Drinks Data Set* consists of a set of records where each entity is an ingredient in a popular bartending recipe found on the Internet.
3. The *IMDB Data Set* is a collection of casts of actors that participated in movies between the years of 1900 and 1960 extracted from the Internet Movie Database
4. The *Citeseer Data* is a set of co-publication records from the Citeseer online library and index of computer science publications. Since the entities are represented by first initial and last name, a single name might correspond to several people.

## 6.2 Empirical Results

We tested our algorithm in a variety of configurations on the datasets listed in Table 3. The results in Table 4 are reported in terms of the average BDeu score, i.e. the final BDeu score obtained by the network averaged over the number of records in the dataset. The number of edges in the resulting Bayes Nets is reported in Table 5. It is interesting to note that the BDeu scores corresponding to the Bayes Nets obtained by running *SBN* as described in Table 1 are very close to the ones obtained by random hillclimbing, but have significantly lower number of edges. This supports our claim that the frequent itemsets indeed contain information most relevant to the construction of the highest scoring Bayes Net. It is evident from the results that each of the proposed augmenting algorithms increase the score. We note however that after augmenting the network with highest-mutual-information edges the total number of arcs almost doubled with the highest relative improvement in score when compared to other proposed augmenting techniques. The hillclimbing seems to improve the score even further though the number of edges is almost quadrupled compared to *SBNS*.

The final score of the DAG produced by *SBNS* depends on user-defined support and maximum Frequent Set size. We have noticed that for Citeseer, IMDB and Institute datasets lowering support and increasing maximum Frequent Set size results in higher BDeu scores. Figure 1 shows score fluctuations when varying maximum Frequent Set size given fixed support for the Citeseer dataset.

### 6.2.1 Maximum Frequent Set Size

In our experiments we tried different maximum Frequent Set sizes: ( $mfss = 2 \dots 5$ ). The lower bound  $mfss = 2$  means that we consider only pairs of items and thus the structure

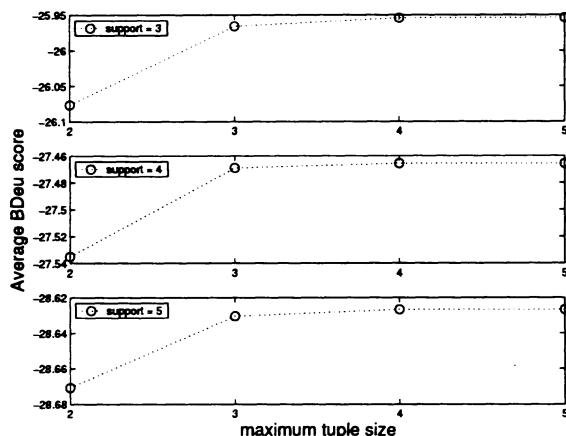


Figure 1: BDeu scores for Citeseer dataset for different parameterizations of the *SBNS* algorithm

Table 4: Average BDeu scores. ( $s = 4$ ,  $mfs = 4$ ; 250,000 random edges considered for hillclimbing)

dataset	rand hlclmb	<i>SBNS</i>	<i>SBNS</i> +MIe	<i>SBNS</i> +MIe+2 <sup>nd</sup>	<i>SBNS</i> +MIe+2 <sup>nd</sup> +hlclmb
citeseer	-33.26	-27.466	-27.375	-27.273	-26.962
imdb	-121.00	-113.15	-112.45	-112.18	-111.28
institute	-11.87	-13.28	-13.18	-13.13	-12.08
drinks	-6.72	-7.21	-7.02	-7.01	-6.705

Table 5: Number of links in the resulting nets. ( $s = 4$ ,  $mfs = 4$ ; 100,000 random edges considered for hillclimbing)

dataset	rand hlclmb	<i>SBNS</i>	<i>SBNS</i> +MIe	<i>SBNS</i> +MIe+2 <sup>nd</sup>	<i>SBNS</i> +MIe+2 <sup>nd</sup> +hlclmb
citeseer	88,259	29,004	48,724	53,790	116,558
imdb	112,773	33,434	52,376	57,236	111,281
institute	1,672	346	398	457	1,159
drinks	723	51	123	133	709

learned is based solely on two-way marginal counts. Figure 1 shows that there is an obvious loss in accuracy when high order interactions are not taken into account. Beyond a maximum Frequent Set size of 4 the number of Frequent Sets does not increase substantially in these datasets and hence the behavior of *SBNS* changes little.

We have to note here, that there is a natural upper bound on the maximum tuple size due to the sparsity of the datasets. For example, there are 94,016 publications in the Citeseer database that have 2 authors and only 3,022 that have exactly 6 authors. The potential

number of publications that have 6 authors, given the total number of authors in the database is  $1.8399e + 27$ , so the empirical number is only  $1.6626e - 22\%$  of the total. The exponential drop in the number of occurrences as the size of the tuples increases is shown on Figure 6.2.1. Hence, we cannot expect a great improvement in the score of the Bayes Net when increasing the maximum tuple size, since there is not enough support for larger tuples.

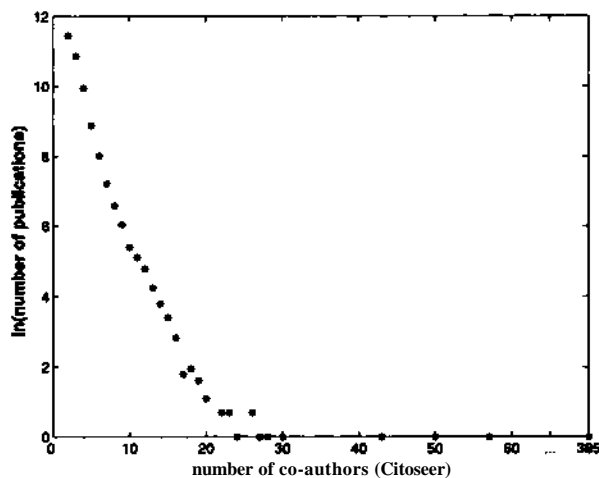


Figure 2: Exponential drop in the number of publications as the number of co-authors increases in the Citeseer Dataset

### 6.2.2 Support

Lowering support greatly increases the number of Frequent Sets to be considered during screening. However, it also introduces quite a few interactions between variables that have low marginal counts. Model fitting in contingency tables in general is sensitive to very low marginal counts even if they are not zero Here we use BDeu, which is less sensitive to low counts. Despite this, it seems to be a good idea to keep support relatively large. In our case, we have tested a few support sizes on smaller datasets and found  $s = 3,4$  to be reasonable support choices. The overall score of the model seems to be better with  $s = 3$ , however it seems to overfit more as is shown in Table 6.

### 6.2.3 Overfitting

We used holdout sets to study overfitting. We withheld roughly a third of the dataset in each case and compared average likelihood per node between the training and testing datasets. The results are summarized in Table 6. The networks learned using *SBNS* always score higher (better) than those learned by hillclimbing on the testing dataset. This indicates that *SBNS* algorithm learns better fitting models. As can be seen from Table 6, the difference in average loglikelihood score for training and testing is in general smaller for hillclimbing. Also, the average loglikelihood of the testing set is worse than the training sets, indicating some degree of overfitting. We believe that some overfitting occurs due to the multiple hypothesis

testing of hundreds of thousands of possible parents. Correction for multiple hypothesis testing problem (similar to corrections used in other learning algorithms such as (Oates & Jensen, 1998)) will be incorporated into *SBNS* in the future.

Table 6: Overfitting testing

dataset	train	test
citeseer hillclimb	-30.6738	-31.0127
citeseer $s = 3$	-23.9227	-26.3253
citeseer $s = 4$	-24.1959	-25.0119
imdb hillclimb	-112.81	-114.851
imdb $s = 3$	-98.1607	-110.499
imdb $s = 4$	-100.203	-107.035

### 6.2.4 Underlying Frequency Distribution

We also note that the algorithm has not performed as well on the Drinks dataset. The dataset seems to correspond to our sparseness requirements, but there is a major difference in the frequency distribution between Drinks and other datasets. A plot of frequency distributions is shown on figure 6.2.4.

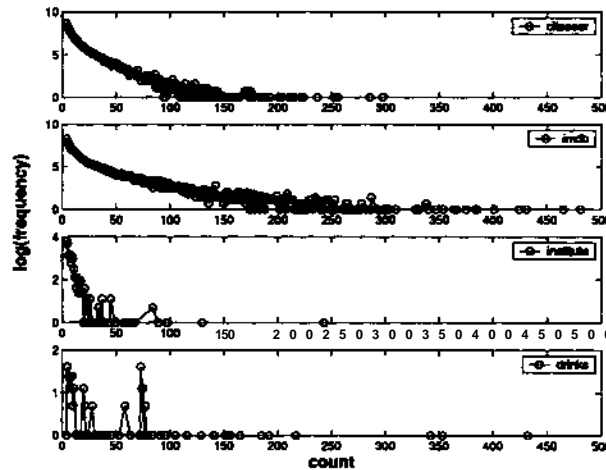


Figure 3:  $\log(\text{Frequency})$  distributions for all datasets, where *count* is the number of times an entity had appeared in the dataset and *frequency* is number of entities that have appeared *count* number of times

It can be observed from the plot that Citeseer, IMDB and Institute datasets that come from social network type sources exhibit close to power law distributions, whereas Drinks has a roughly linear dependency between the number of connections per entity and the

Table 7: Total times for random hillclimbing, *SBNS* and *SBNS+UIe* to create a Bayes Net (in mins). ( $s = 4$   $mfss = 4$ )

dataset	rand hillclmb	<i>SBNS</i>	<i>SBNS+MIe</i>
citeseer	171	59.8	87
imdb	193	225.6	252.8
institute	.53	.016	.017
drinks	.37	.016	.017

number of entities with the same number of connections (frequency). It is thus far only a hypothesis, that our algorithm will show better performance if the dataset were close to power law distribution.

### 6.2.5 Time performance

All experiments were conducted on unloaded 2GHz Pentium IV machines with 2GB of RAM. The total times required to run the algorithm and the time it took random hillclimbing to create a Bayes Net by adding/removing/reversing 250,000 edges are reported in Table 7. We also break the total time into segments corresponding to major steps of the algorithm as reported in Table 8.

Table 8: Time (min) per task for *SBNS*. ( $s = 4$ ,  $mfss = 4$ )

dataset/task	freq sets	lcl strct search	<i>Edwnp</i> & DAG	MI augment	2 <sup>nd</sup> degree augment
citeseer	65.49	4.11	.2	17.2	97.5
imdb	196.22	15.43	13.93	27.23	22.33
institute	.00	.02	.00	.00	.02
drinks	.00	.00	.01	.00	.00

The biggest cost is to obtain the frequencies; the time it takes to perform the remaining operations depends on the number of Frequent Sets that occur more frequently than pre-defined support. Our experiments have shown that number to be only a small fraction of the total number of entities (nodes). It is also interesting to note that random hillclimbing is very fast while the network consists of many small subgraphs, but as soon as the subgraphs are joined together by new edges, the time increases tremendously due to the complexity of cycle detection. For example, it takes random hillclimbing on the order of 10 minutes to add/remove/reverse 250,000 edges, but it takes over 6 hours to perform those operations given the same number of nodes for 300,000 edges with relatively small increase in the score. In that sense, the random graphs might not be exactly random as discussed in (Callaway et al., 2001).

## 6.3 Application

One of the important and growing application fields of large Bayes Nets is recommender systems. The purpose of the service is to provide user with suggestion of products that

he/she is likely to buy based on their historical preferences. One of the well known and very successful recommender systems is Amazon. We did not have Amazon data available to us, but we simulated the query based on the Citeseer dataset. The mapping is as follows: suppose that the set of co-authors of a paper represents user's preferences of particular products. We then learn a Bayes Net based on the available co-authorship information and query the network with incomplete subsets of authors to predict the most likely selection of entities (or authors in our case) that completes the given set.

To answer the query we simply calculated the loglikelihoods of the most likely completions and reported the top  $n$  with the highest scores. The set of most likely completions are formed from "similar" and "popular" entities. Entities are considered similar if they are in close proximity to the query items in the network and the popular items are the ones with high marginal counts. It is particularly reasonable to make those assumptions in the recommender system case, as people are more likely to choose something that is close to what they like or something popular than any randomly picked item.

Here is an example of a query and its completion. The query is a subset of former or present members of Daphne Roller's group (DAG):  $\{d\ koller, I\ getoor, a\ pfeffer, b\ taskar\}$ . Results are presented in Tables 9 and 10.

Table 9: 3 most likely completions of size 1 for 4 members of the DAG group

completion	score
koller <b>friedman</b> pfeffer getoor taskar	-22.523647
koller pfeffer getoor <b>tong</b> taskar	-22.694517
koller <b>boyen</b> pfeffer getoor taskar	-23.079099

Table 10: 3 most likely completions of size 1 for 4 members of the DAG group

completion	score
koller <b>grove halpern</b> pfeffer getoor taskar	-24.065985
koller <b>malik weber</b> pfeffer getoor taskar	-24.335174
koller <b>russell parr</b> pfeffer getoor taskar	-24.688802

The suggested completions are in fact people that are either part of or close collaborators of Daphne Roller's group, thus by analogy we might expect a set of relevant items to be predicted by the recommender system using this algorithm. It is interesting to note that in the example above the one most likely person to complete the given subset is different (Table 9) than the suggestions provided by the algorithm under the assumption of 2 missing people (Table 10). This observation suggests that there are more complex interactions that could not be found by systems built on pairwise statistics. The inference took less than a second.



## 7 Related Work

An increase in the amounts of collected data has facilitated an interest in modelling massive datasets. There are several approaches to dealing with large data collections. The most relevant to our work is the direction of fast learning from sparse data.

Some of the earlier work in this area has concentrated on efficient representation of sparse data and caching of n-way counts (Moore & Lee, 1998). (Chickering & Heckerman, 1999) and (Meila, 1999) have noted that computations requiring one-way and pairwise counts can be sped up significantly when dealing with sparse data using caching and such data structures as ADTrees (Moore & Lee, 1998). We believe that this body of work has great potential and thus we build on the ideas introduced in these papers by utilizing the sparse data representation and low overhead efficient calculation of the marginals.

Using frequent sets when learning Bayes Nets on the local scale was also explored in (Pavlov et al., 2003). The goal of this work was to answer probabilistic queries on a subset of variables, thus there was no need to combine local information to obtain the joint distribution once the query size was estimated. The performance of Bayes Nets learned from a selection of variables was reported to be worse though close in accuracy to the inferences drawn from a Bayes Net learned on a full dataset. In (Hollmen et al., 2003) it has been proposed to integrate Frequent Sets as a local methodology when modelling joint distributions. This work has shown that mixture models obtained from Frequent Sets using maximum entropy are more accurate, thus supporting our claim that frequent sets contain important local information when modelling joint distributions.

One approach to speed up structural search in Bayes Nets for massive datasets has been to restrict the possible number of parents. The full Sparse Candidate Algorithm is presented in (Friedman et al., 1999). In its original form it is a method to speed up hillclimbing at the cost of lower performance, though in practice the performance loss was shown to be not so great. This work is yet another motivation for us, since structural search on the local scale inadvertently restricts the number of parents. However, since on the global scale the number of parents in our Bayesian Network is not limited we perceive it as an improvement on the original Sparse Candidate algorithm.

Sampling was proposed as one of the techniques to speed up modelling in massive datasets in (Kaebling, 1990; Maron & Moore, 1997; Hulten & Domingos, 2002; Pelleg & Moore, 2002). Though an interesting direction it seems to be orthogonal to our approach.

The idea of augmenting Bayes Nets with high mutual information edges between entities is based on the fact that such dependencies could not be accounted for in frequent sets. The fast computation used in this work is based on (Meila, 1999).

Since the proposed algorithm does not obtain an optimal network, it could be used as a structure initialization for other Bayesian Network learning algorithms as described in (Cooper & Herskovits, 1991; Buntine, 1991; Spiegelhalter et al., 1993; Heckerman et al., 1995; Moore & Wong, 2003) and other algorithms.

## 8 Conclusion

We have presented a tractable solution to the Bayes Net structure search problem in sparse datasets. Like other researchers, we use Frequent Sets to take advantage of sparseness. Our main new contribution is to perform structural search on the local level in order to produce the global model. We propose several techniques to improve the score of the created net. One of the key improvements is augmentation by edges with high mutual information for entities that have not co-occurred in the dataset.

We have performed an empirical study of *SBNS* using two small and two large (over  $10^5$  attributes) datasets. We show tractable times while maintaining accuracy better than hillclimbing, which is the only tractable alternative for learning structure in networks of this size. Empirical evidence also shows that higher accuracies are achieved without requiring more complex structures. It seems likely that when it comes to using the Bayes Net for inference as suggested in Section 1, the relatively small number of edges in such networks will be advantageous in comparison with networks obtained from hillclimbing.

We believe that SBNS serves two primary purposes. First, it opens new horizons for modelling joint distributions of massive transactional datasets. Second, it can be viewed as a novel way to postprocess Frequent Sets in commercial data mining.

Third, we raise the question of structural search that takes into consideration characteristics of the dataset being modelled. Model selection is greatly effected by properties such as frequency distribution. We believe, that there is immense potential in exploiting those properties to obtain high accuracy models in a fraction of time required for generic techniques.

## 9 Appendix

Here we provide two illustrations to support our claim that the Frequent Sets contain essential information needed to build a Bayes Net from sparse data. First, we show that in sparse large datasets positive correlation between two variables is much stronger than negative. The second shows the advantage of considering Frequent Sets of order higher than 2.

### 9.1 Frequent Sets are useful

Suppose we have 2 binary variables  $x$  and  $y$ . Assume our dataset is sparse and has  $R$  records, where  $R$  is very large.

Let us look at the correlation coefficient  $\rho$  of the two variables. Under the multinomial sampling model the observed correlation coefficient

$$r = \frac{N_{\bar{x}\bar{y}}N_{xy} - N_{\bar{x}y}N_{x\bar{y}}}{\sqrt{N_{\bar{x}+}N_{x+}N_{+\bar{y}}N_{+y}}} \quad (4)$$

is the maximum likelihood estimate of  $\rho$  (Bishop et al., 1977).

Case 1: Two entities have co-occurred with each other  $v$  times and  $kv$  times separately elsewhere in the dataset, where  $k \rightarrow 0$ . Then

$$\begin{aligned}
 r &= \frac{v(R - 2kv - v) - (kv)^2}{\sqrt{(v + kv)^2(R - kv - v)^2}} \\
 &= \frac{J}{1 + kR - kv}, \text{ as } k \rightarrow 0 \tag{5}
 \end{aligned}$$

In fact, only as  $A; \rightarrow \sqrt{v}/f \sim 1$  (which is clearly a violation of the sparseness assumption), the correlation between  $x$  and  $y$  becomes significant.

Case 2: Two entities have occurred with frequency  $Kv$  but never with each other.  $K$  in this case could be rather large, but still conforming to the sparseness assumption, i.e.  $Kv < C R_l$  then  $r = -\frac{1}{K} \approx 0$ . Note that when  $K = k + 1$ , we have the same frequency of occurrence as in Case 1, yet only if  $K \rightarrow \frac{C}{v}$  the correlation would become significant.

This means that assuming sparseness, positive correlations are much stronger than negative ones. We have obtained the same results comparing BDeu scores for  $x \rightarrow y$  models for the above 2 cases. Thus, when learning a Bayes Net we are much more likely to increase the score by screening Frequent Sets first.

## 9.2 Benefit of higher order Frequent Sets

Suppose we have 3 variables  $A$ ,  $B$  and  $C$ . Let their counts be  $N_A = 971$ ,  $N_B = 936$ ,  $N_C = 156$ ,  $N_{AB} = 96$ ,  $N_{AC} = 97$ ,  $N_{BC} = 51$ ,  $N_{ABC} = 25$ . Then, if we only consider pairs of items, we will find that  $A$  and  $B$  as well as  $A$  and  $C$  are independent and model  $B \rightarrow C$  is the best DAG for the pair  $BC$  according to the BDeu score. This results in DAG on Figure 4(a). However, when we look at the triplet, we find that the model on Figure 4(b) is the one that fits the data best. This shows that a model that considers only two-way counts may be inaccurate. Similar analysis can be done for higher order interactions.

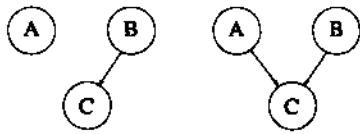


Figure 4: Best fitting models when considering (a) only pairs of items, and (b) triples.

## References

- Agrawal, R., Imielinski, T., & Swami, A. (1993). Mining association rules between sets of items in large databases. *ACM SIGMOD 12* (pp. 207-216).
- Agrawal, R., & Srikant, R. (1994). Fast algorithms for mining association rules. *VLDB 20* (pp. 487-499).

- Barabasi, A. (2002). *Linked: The new science of networks*. Perseus Publishing.
- Bishop, Y., Fienberg, S., & Holland, P. (1977). *Discrete multivariate analysis: Theory and practice*. MIT Press.
- Breese, J., Heckerman, D., & Kadie, C. (1998). Empirical analysis of predictive algorithms for collaborative filtering. *UAI14*.
- Breiger, R. (2003). Emergent themes in social network analysis: Results, challenges, opportunities. *Dynamic Social Network Modeling and Analysis: Workshop Summary and Papers*.
- Buntine, W. (1991). Theory refinement on Bayesian networks. *UAI 7* (pp. 52–60).
- Callaway, D., Hopcroft, J., Kleinberg, J., Newman, M., & Strogatz, S. (2001). Are randomly grown graphs really random? *Physical Review*.
- Chickering, D., & Heckerman, D. (1999). Fast learning from sparse data. *UAI 15*.
- Cooper, G., & Herskovits, E. (1991). A Bayesian method for constructing Bayesian belief network from databases. *UAI 7* (pp. 86–94).
- Friedman, N. (2004). Inferring cellular networks using probabilistic graphical models. *Science*.
- Friedman, N., Nachman, I., & Pe'er, D. (1999). Learning bayes network structure from massive datasets: The "sparse candidate" algorithm. *UAI 15* (p. 206:215).
- Han, J., & Kamber, M. (2000). *Data mining: Concepts and techniques*. Morgan Kaufmann Publishers.
- Heckerman, D., Geiger, D., & Chickering, D. (1995). Learning Bayesian Networks: The combination of knowledge and statistical data. *Machine Learning, 20*, 197–243.
- Hollmen, J., Seppanen, J., & Mannila, H. (2003). Mixture models and frequent sets: combining global and local methods for 0-1 data. *SIAM ICDM*.
- Hulten, G., & Domingos, P. (2002). Mining complex models from arbitrarily large databases in constant time. *ACM SIGKDD 8*.
- Kaebling, L. (1990). *Learning in embedded systems*. Doctoral dissertation, Stanford University.
- Mannila, H., & Toivonen, H. (1996). Multiple uses of frequent sets and condensed representations. *KDD 2* (pp. 189 – 194).
- Maron, O., & Moore, A. (1997). The racing algorithm: Model selection for lazy learners. *Artificial Intelligence Review* (pp. 193–225).

- Meila, M. (1999). *An accelerated Chow and Liu algorithm: fitting tree distributions to high dimensional sparse data* (Technical Report AIM-1652). MIT.
- Moore, A., & Lee, M. S. (1998). Cached sufficient statistics for efficient machine learning with large datasets. *Journal of Artificial Intelligence Research*, 8, 67-91.
- Moore, A., & Wong, W. (2003). Optimal Reinsertion: A new search operator for accelerated and more accurate Bayesian network structure learning. *ICML 20*.
- Moreno, J., & Jennings, H. (1938). Statistics of social configuration. *Sociometry*, 342-374.
- Oates, T., & Jensen, D. (1998). Large datasets lead to overly complex models: An explanation and a solution. *KDD 4*.
- Pavlov, D., Mannila, H., & Smyth, P. (2003). Beyond independence: probabilistic models for query approximation on binary transaction data. *IEEE Transactions on Knowledge and Data Engineering*.
- Pelleg, D., & Moore, A. (2002). Using taxjan's red rule for fast dependency tree construction. *NIPS 15*.
- Spiegelhalter, D., Dawid, A., Lauritzen, S., & Cowell, R. (1993). Bayesian analysis in expert systems. *Statistical Science*, 219-282.