

Learning Bayesian Networks with Low Rank Conditional Probability Tables

Adarsh Barik, Jean Honorio

abarik@purdue.edu, jhonorio@purdue.edu

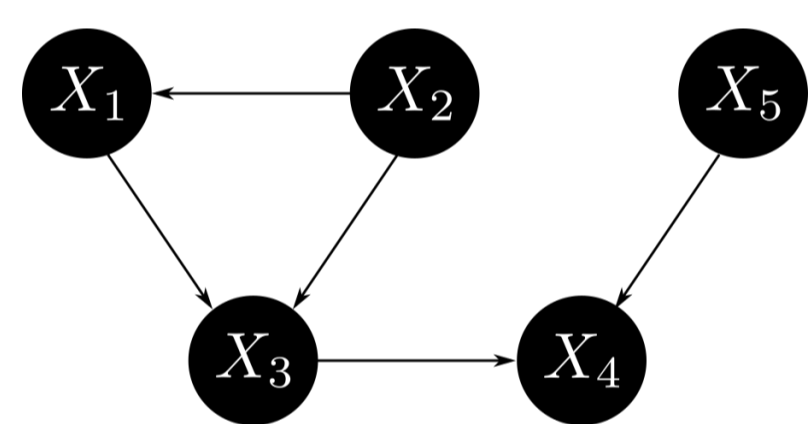
Department of Computer Science, Purdue University, West Lafayette - IN, 47907.

OVERVIEW

- Bayesian networks can be used to model complex interactions amongst constituent variables of a system
- Question:** Can we learn the structure (directed) of a Bayesian network by observing variables of a system?
- Answer:** NP-hard in general. We propose a method to learn exact structure of a class of Bayesian networks by making black-box queries.
- Theoretical Guarantees:** We formally prove the correctness of our method and provide polynomial bounds on sample and time complexity.

PRELIMINARIES

- Bayesian network:** A Bayesian network represents a joint probability distribution over the set of random variables defined on the nodes of DAG which factorizes according to the DAG structure, i.e., $P(X_1, X_2, X_3, X_4, X_5) = P(X_1|X_2)P(X_2)P(X_3|X_1, X_2)P(X_4|X_3, X_5)P(X_5)$.



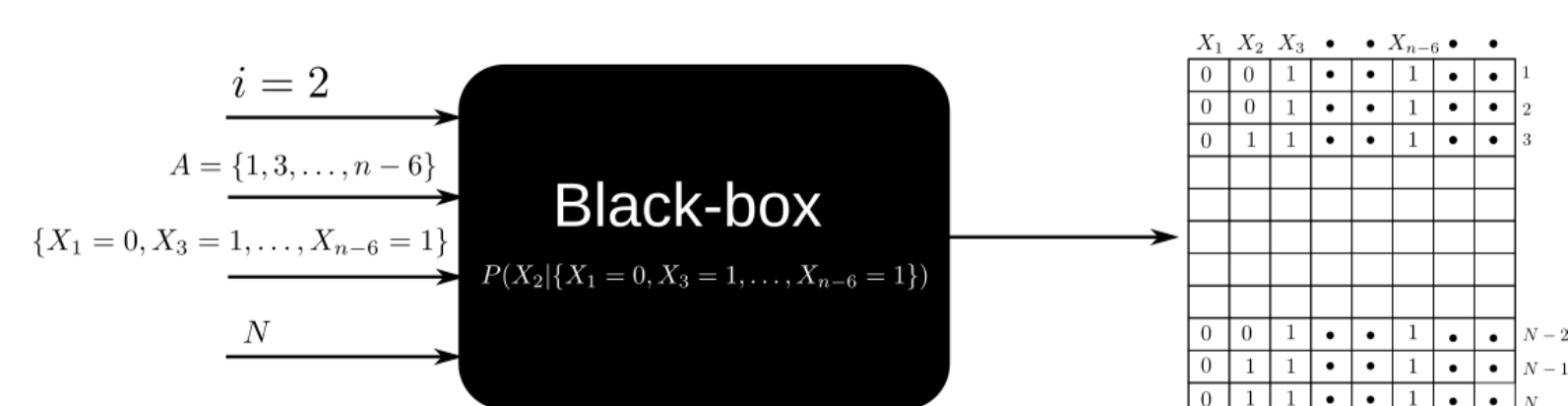
- Rank k conditional probability distribution** A node i of a Bayesian network is said to be rank k representable with respect to a set $A \in V \setminus \{i\}$ and probability distribution P if $\forall x_i \in \text{dom}(X_i), x_{A(i)} \in \text{dom}(X_{A(i)})$,

$$P(X_i = x_i | X_A = x_A) = \sum_{\substack{S \subseteq \{i\} \cup A \\ 1 \leq |S| \leq k, i \in S}} Q_S(X_S = x_S)$$

		$X_1 X_2$				X_1			X_2	
		00	01	10	11	0	1	0	1	
X_3	0	$\frac{13}{36}$	$\frac{37}{90}$	$\frac{5}{18}$	$\frac{59}{180}$	$\frac{1}{9}$	$\frac{1}{6}$	$\frac{1}{12}$	$\frac{1}{12}$	$\frac{2}{15}$
	1	$\frac{23}{36}$	$\frac{53}{90}$	$\frac{13}{18}$	$\frac{121}{180}$	$\frac{2}{9}$	$\frac{1}{6}$	$\frac{1}{4}$	$\frac{1}{4}$	$\frac{1}{5}$
		$P(X_3 X_1, X_2)$				$Q_3(X_3)$	$Q_{31}(X_3, X_1)$		$Q_{32}(X_3, X_2)$	

Example of Rank-2 CPT

PROBLEM DESCRIPTION

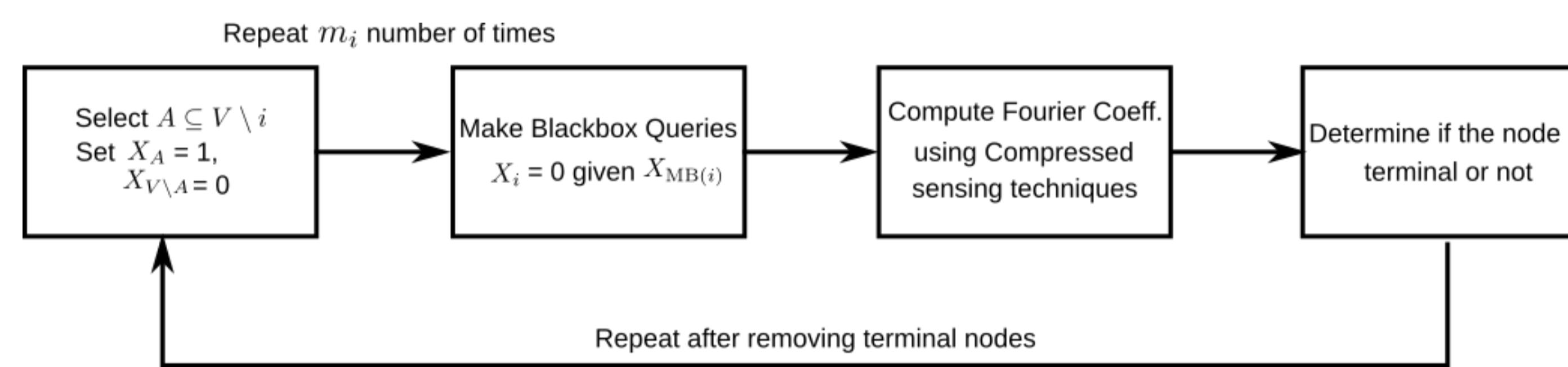


Can we recover the DAG structure using *blackbox queries* with theoretical guarantees of correctness and efficiency in terms of time and sample complexity?

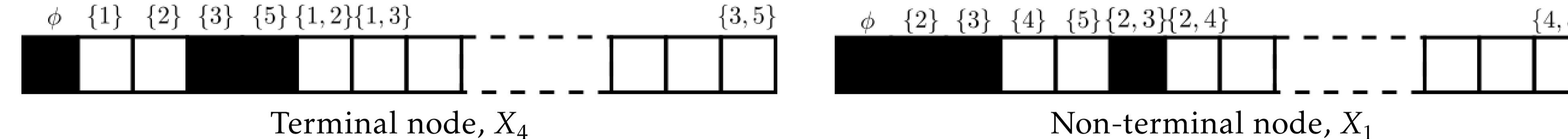
OUR CONTRIBUTION

- Introduction of Low rank CPTs:** CPT can be treated as summation of multiple simple tables, each of them depending only on a handful of parents
 - Connect this notion of rank of a CPT to the Fourier transformation of a specific real valued set function [1]
- For a set function $f : 2^T \rightarrow \mathbb{R}$, $f(A) = \sum_{B \in 2^T} \hat{f}(B) (-1)^{|A \cap B|}$ where Fourier coefficient $\hat{f}(B) = 2^{-|T|} \sum_{A \in 2^T} f(A) (-1)^{|A \cap B|}$.
- Use compressed sensing techniques to show that the Fourier coefficients of this set function can be used to learn the structure of the Bayesian network

Algorithm



Differentiate between Terminal and Non-terminal node



MAIN THEOREM

Let $g_i(A_j)$ be an approximation for $f_i(A_j) = P(X_i = 0 | X_{MB(i)})$ by taking only upto pairwise Fourier coefficients. Let $|f_i(A_j) - g_i(A_j)| \leq \epsilon_j, \forall A_j \in \mathcal{A}_i$ for some $\epsilon_j > 0$. \hat{g}_i by solving the following optimization problem for each node i .

$$\beta_i = \min_{\hat{g}_i \in \mathbb{R}^{|\rho_i|}} \|\hat{g}_i\|_1 \quad \text{s.t.} \quad \|\mathcal{M}_i \hat{g}_i - f_i\|_2 \leq \epsilon \quad \text{where} \quad \epsilon = \sqrt{\sum_{A_j \in \mathcal{A}_i} \epsilon_j^2}. \quad (1)$$

Let \mathcal{A}_i be a collection of m_i sets $A_j \in 2^{V-i}$ chosen uniformly at random and $\mathbf{g}_i \in \mathbb{R}^{m_i}$ be a vector whose j^{th} row is $g_i(A_j)$ and $\hat{\mathbf{g}}_i \in \mathbb{R}^{n+\binom{n-1}{2}}$ be a vector with elements of form $\hat{f}_i(B_k) \forall B_k \in \rho_i$ where $\rho_i = \{B_k \mid B_k \in 2^{\bar{V}}, |B_k| \leq 2\}$ is a set which contains support(\hat{f}_i). Consider, $\mathbf{g}_i = \mathcal{M}_i \hat{\mathbf{g}}_i$ where, $\mathcal{M}_i \in \{-1, 1\}^{m_i \times n}$ such that $\mathcal{M}_{jk}^i = (-1)^{|A_j \cap B_k|}$.

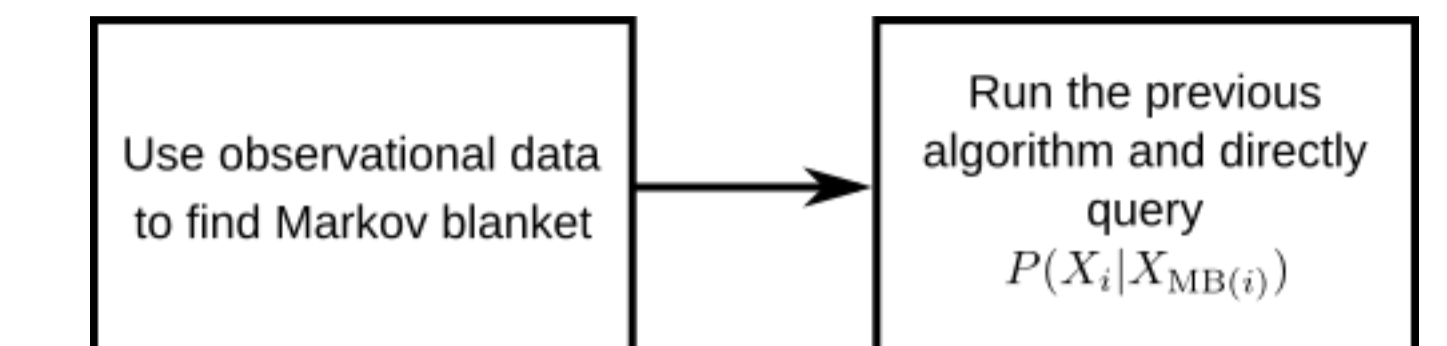
Theorem 1. Suppose $\hat{\mathbf{g}}_i$ is constructed by computing $\hat{g}_i(B_k)$ using B_k from a fixed collection ρ_i and \mathbf{g}_i is computed by selecting m_i sets A_j uniformly at random from $2^{\bar{V}}$. Let $m_i \geq \mathcal{O}(\max(|\text{support}(\hat{\mathbf{g}}_i)| \log^4(n + \binom{n-1}{2}), |\text{support}(\hat{\mathbf{g}}_i)| \log \frac{1}{\delta}))$ and β_i be solved using equation (1). Then with probability at least $1 - \delta$, we have $\|\beta_i - \hat{\mathbf{g}}_i\|_2 = \mathcal{O}(\frac{\epsilon}{\sqrt{m_i}})$. If the minimum non-zero element of $|\hat{\mathbf{g}}_i|$ is $\Omega(\frac{\epsilon}{\sqrt{m_i}})$ then β_i recovers $\hat{\mathbf{g}}_i$ up to the signs. Furthermore, if non-terminal nodes are not rank 2 then $|\beta_i(B)| = \mathcal{O}(\frac{\epsilon}{\sqrt{m_i}}), \forall B \in \rho_i, |B| = 2$ iff i is a terminal node and $\hat{\pi}(i) = \{B \mid |B| = 1, |\beta_i(B)| = \Omega(\frac{\epsilon}{\sqrt{m_i}})\}$ correctly recovers the parents of node i .

RESULTS

Algorithms	Sample Complexity	Time Complexity	Selections	Queries
Our Work (no observational data)	Blackbox - $\mathcal{O}(nk^3 \log^4 n (\log k + \log \log n))$	$\mathcal{O}(n^4 k \sqrt{n} \log n)$	$\mathcal{O}(n)$	$\mathcal{O}(nk^3 \log^4 n)$
Our Work (with observational data)	Observational - $\mathcal{O}(n)$, Blackbox - $\mathcal{O}(nk^3 \log^5 k)$	$\mathcal{O}(n^4), \mathcal{O}(nk^4 \sqrt{k} \log k)$	$\mathcal{O}(n)$	$\mathcal{O}(nk^3 \log^4 k)$
Bello et. al, 2018	Interventional - $\mathcal{O}(n^2 2^k \log n)$	$\mathcal{O}(n^2 2^k \log n)$	$\mathcal{O}(n^2)$	$\mathcal{O}(n^2 2^k)$
Kocaoglu et. al, 2017	Interventional - no guarantees	$\mathcal{O}(2^n kn^2 \log^2 n)$	$\mathcal{O}(\log n)$	$\mathcal{O}(2^n \log n)$

Here n is the number of nodes, k is the maximum size of the Markov blanket. The maximum number of parents of a node is $\mathcal{O}(k)$.

USING OBSERVATIONAL DATA

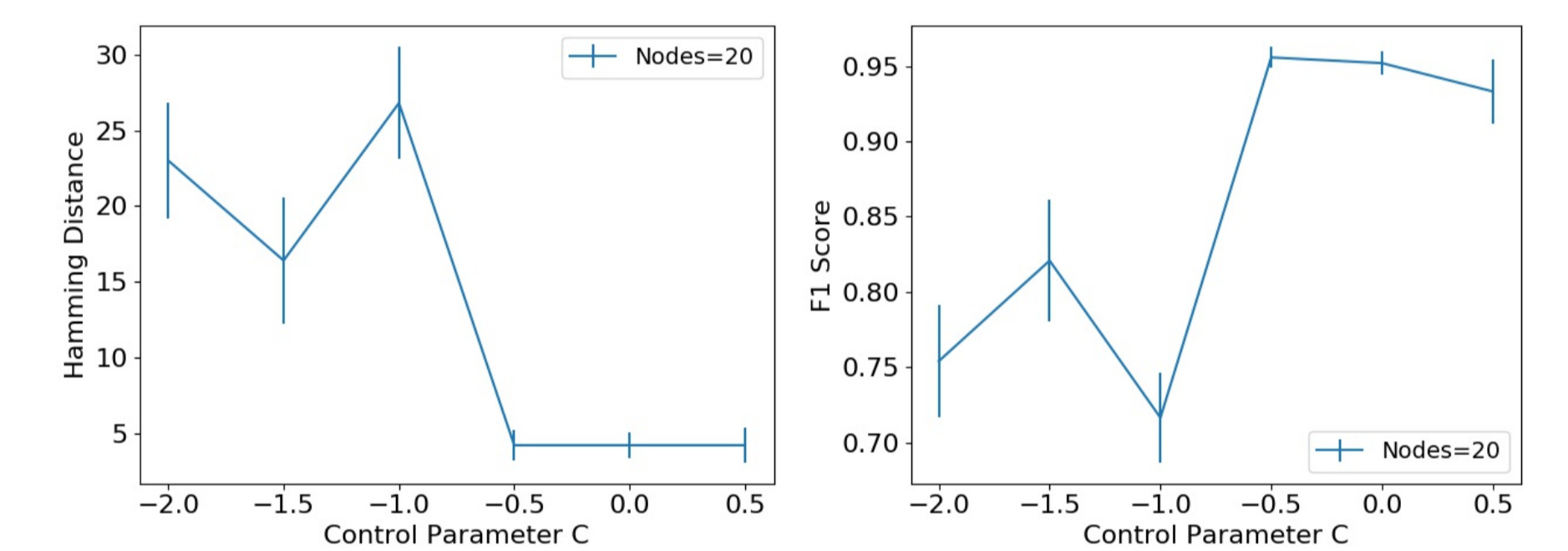


Let $\mathbb{P} = \{P \text{ is faithful to } G, |P(X_i|X_l) = P^*(X_i|X_l), \forall i, l \in \{1, \dots, n\}\}$. If there exists a probability distribution $\hat{P} \in \mathbb{P}$ such that each node i is rank 2 with respect to $\text{MB}_G(i)$ and \hat{P} , then we can recover the Markov blanket by solving the following system of equations:

$$P(X_i = 0, X_j = 0) = \bar{Q}_i(X_i = 0)P(X_j = 0) + \sum_{\substack{j \in \bar{V} \\ j \neq i}} \bar{Q}_{ij}(X_i = 0, X_j = 0)P(X_j = 0, X_i = 0) + \bar{Q}_i(X_i = 0, X_j = 0)P(X_j = 0) \forall l = \{1, \dots, n\}, l \neq i$$

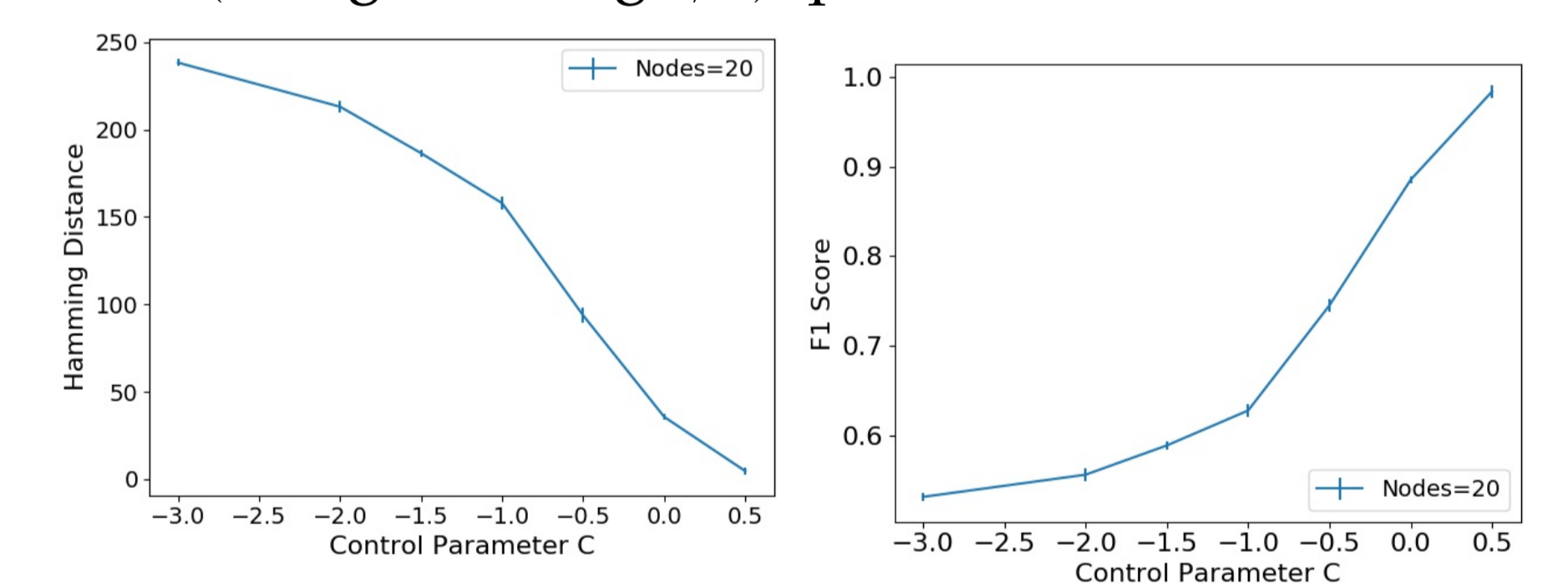
$$P(X_j = 0) = \bar{Q}_j(X_j = 0) + \sum_{\substack{j \in \bar{V} \\ j \neq l}} \bar{Q}_{lj}(X_l = 0, X_j = 0)P(X_j = 0)$$

VERIFICATION



(a) Hamming distance with control parameter C (b) F1 score with control parameter C

Figure 1: Without observational data. $m_i = 10^C \max(k^2 \log^4 n', k^2 \log 1/\delta)$ queries for each node i .



(a) Hamming distance of Markov blanket recovery with control parameter C (b) F1 score of Markov blanket recovery with control parameter C

Figure 2: With observational data. $N = \max(10^{\frac{C \log n}{e^2}}, n)$ observational samples

References

- [1] Peter Stobbe and Andreas Krause. Learning Fourier Sparse Set Functions. In *Artificial Intelligence and Statistics*, pages 1125–1133, 2012.