

Estimating Causal Effects from Learned Causal Networks



Anna K Raichev(araichev@uci.edu), Jin Tian(itian@iastate.edu),
Alex Ihler(ihler@ics.uci.edu), and Rina Dechter (dechter@ics.uci.edu)



Overview

We propose an alternative to the estimand based paradigm for answering causal queries. The idea is to learn the full causal model from the observational data and causal diagram, and then answer the query by applying Probabilistic Graphical Models (PGM) algorithms. We show that this model completion learning approach can be far more effective than estimated approaches, particularly in large models when the estimand computation is complex and the induced width of the diagram is small.

Contributions:

1. Provide a first of its kind, extensive empirical evaluation on causal effect algorithms on varied and large synthetic and real networks.
2. Show empirically that our approach has more accurate estimates than estimated based schemes.

Problem

Given a causal diagram, an identifiable query $P(Y | do(X = x))$ and samples from the observed distribution, the task is to output the distribution of $P(Y | do(X = x))$.

Current Practice

1. Apply state of the art algorithms for identifiability. These are polynomial algorithms involving the graph and the query only. [Tian, 2002]
2. Generate an estimand, namely an algebraic expression for the query involving only probabilistic expressions over the visible variables.
3. Estimate the estimand from the observational data.

Limitations

1. More sophisticated statistical estimation techniques don't scale when functions in the estimand are too large.
2. We can use the *Plug-In method*, in which each term is estimated only on the configurations seen in the observed data. However, this approach also limits the quality of our estimates.

Motivating Example

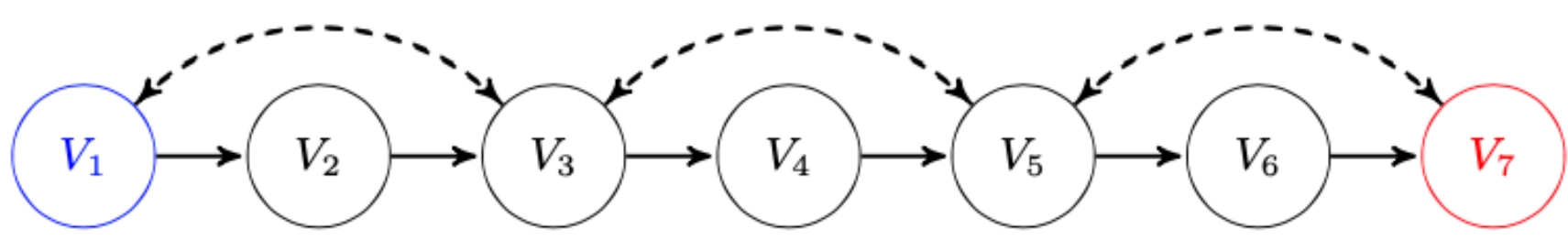


Figure 1: Chain Model with 7 observable variables and 3 latent variables

Using the estimated based approach we get the expression:

$$P(V_7 | do(V_1)) = \sum_{V_2, V_3, V_4, V_5, V_6} P(V_6 | V_1, V_2, V_3, V_4, V_5) P(V_4 | V_1, V_2, V_3) P(V_2 | V_1) \times \prod_{V_i} P(V_i | V_1, V_2, V_3, V_4, V_5, V_6) P(V_5 | V_1, V_2, V_3, V_4) P(V_3 | V_1, V_2) P(V_1)$$

- As model size increases, we have scalability issues.
- However, the induced width of this model is only 3.

Background

Structural Causal Model: $M = \langle U, V, F, P(U) \rangle$

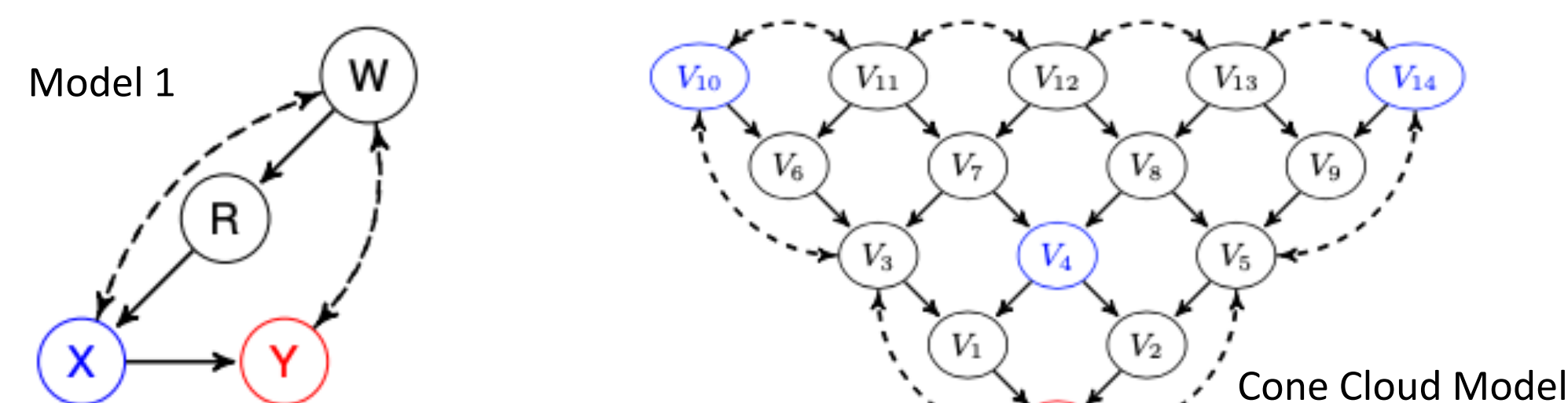
- $U = \{U_1, \dots, U_k\}$ set of unmeasurable latent variables
- $V = \{V_1, V_2, \dots, V_n\}$ set of observable variables
- $F = \{f_i : V_i \in V\}$ is a set of functional mechanisms f_i that each determine the value v_i of their corresponding V_i as a function of V_i 's causal parents $PA_i \subseteq U \cup V \setminus V_i$
- $P(U)$ is a probability distribution over the exogenous variables

Causal Diagram: A SCM M can be associated with a directed graph $G = \langle V \cup U, E \rangle$ called a causal diagram. Each node in the graph uniquely corresponds to a variable in the SCM. There is an arc from node $X \in (U \cup V)$ to node $V_i \in V$ iff $X \in PA_i$

Causal effect and the truncation formula: We use $P(Y | do(X))$ to denote the distributions resulting from an intervention which fixes the value of X , and is called the causal effect of $do(X)$ on Y

$$P(V, U | do(X)) = \prod_{V_j \in X} P(V_j | PA_j) \cdot P(U)$$

Causal Diagrams:



Blue variables are intervened on and red variables are the outcome variables corresponding to the query $P(Y | do(X))$

Learning for Causal Inference

Identifiability

- Any two models that agree on the observational distribution and causal diagram will also agree on $P(Y | do(X = x))$.

EM for Causal Inference (EM4CI)

Algorithm 1: EM4CI

input : A causal diagram $G = \langle U \cup V, E \rangle$, U latent and V observables; D samples from $P(V)$;

output : Estimated $P(Y | do(X = x))$

// $k =$ latent domain size, $BIC_B = BIC$ score of B , D ,
// LL_B is the log-likelihood of B , D

1. Initialize: $BIC_B \leftarrow \inf$,
2. If \neg identifiable(G, Q), terminate.
3. **for** $k = 2, \dots$, to upper bound, **do**
4. $(B', LL_{B'}) \leftarrow \max_{LL} \{EM(G, D, k) \mid \text{for } i = 1 \text{ to } 10\}$
5. Calculate $BIC_{B'}$ from $LL_{B'}$
6. **if** $BIC_{B'} \leq BIC_B$,
7. $B \leftarrow B', BIC_B \leftarrow BIC_{B'}$
8. **else**, break.
9. **endfor**
- 10: $B_{X=x} \leftarrow$ generate truncated CBN from B .
- 11: **return** \leftarrow evaluate $P_{B_{X=x}}(Y)$

1. Check if query is identifiable.
2. Using samples from the observed distribution $P(V)$ to learn a full causal Bayesian network B with domain size k consistent with $(G, P(V))$ using the the resulting model with maximum log likelihood from running the EM algorithm ran 10 times.
3. Compute the BIC score and increase k .
4. Stop when we find the minimum BIC score.
5. truncate M into the causal model M_x by removing the function associated with X and assigning $X = x$ in all functions where X appears.
6. Apply a PGM algorithm to answer apply a PGM algorithm to answer the associated query $P(Y | X = x)$.
7. Return $P(Y | do(X = x))$.

Complexity

- Time and memory are exponential in the induced width.

Benefits & Challenges

Challenges

1. In order to learn the full model we need to learn a domain size for the latent variables.
2. There exists theoretical bounds on sufficient domain sizes. However the bounds are very conservative & can be very large to be practical [J. Zhang et al, 2022].
3. EM algorithm can be slow and converge to incorrect local optima in high dimensional space.

Benefits

1. Learning phase only needs to be performed once to answer any identifiable of form $P(Y | do(X = x))$; traditionally a new estimand would need to be derived for each query.
2. EM4CI consistently yields extremely accurate results.

Experimental Setup

Benchmarks

- Each benchmark includes a causal diagram, a query, and observational data synthetically generated from the full model.
- Used a range of domain sizes of for the variables.
- Test on bayesian networks from real world domains, and created latent confounders from the source vertices.

Performance Measures

- To evaluate the accuracy of $P(Y | do(X = x))$, we use the mean absolute deviation (*mad*): averaging the absolute error over all single-value queries over all instantiations of the intervened and queried variables.
- BIC score is used to evaluate fitness of the learned model and impose some regularization over the domain sizes.

Notation

- Capital letters (X) represent variables, & small letters (x) represent their values. Boldfaced capital letters (\mathbf{X}) denote a collection of variables.
- $n = |V|$, $d = |D(V)|$, $k = |D(U)|$ in the true model, and $k_{lrm} = |D(U)|$ the latent domain of the learned model

Empirical Analysis

Baseline Comparison: Plug In Method

Table 3: Results of EM4CI & Plug-In on $P(Y | do(X))$ ($d, k = (2, 10)$)

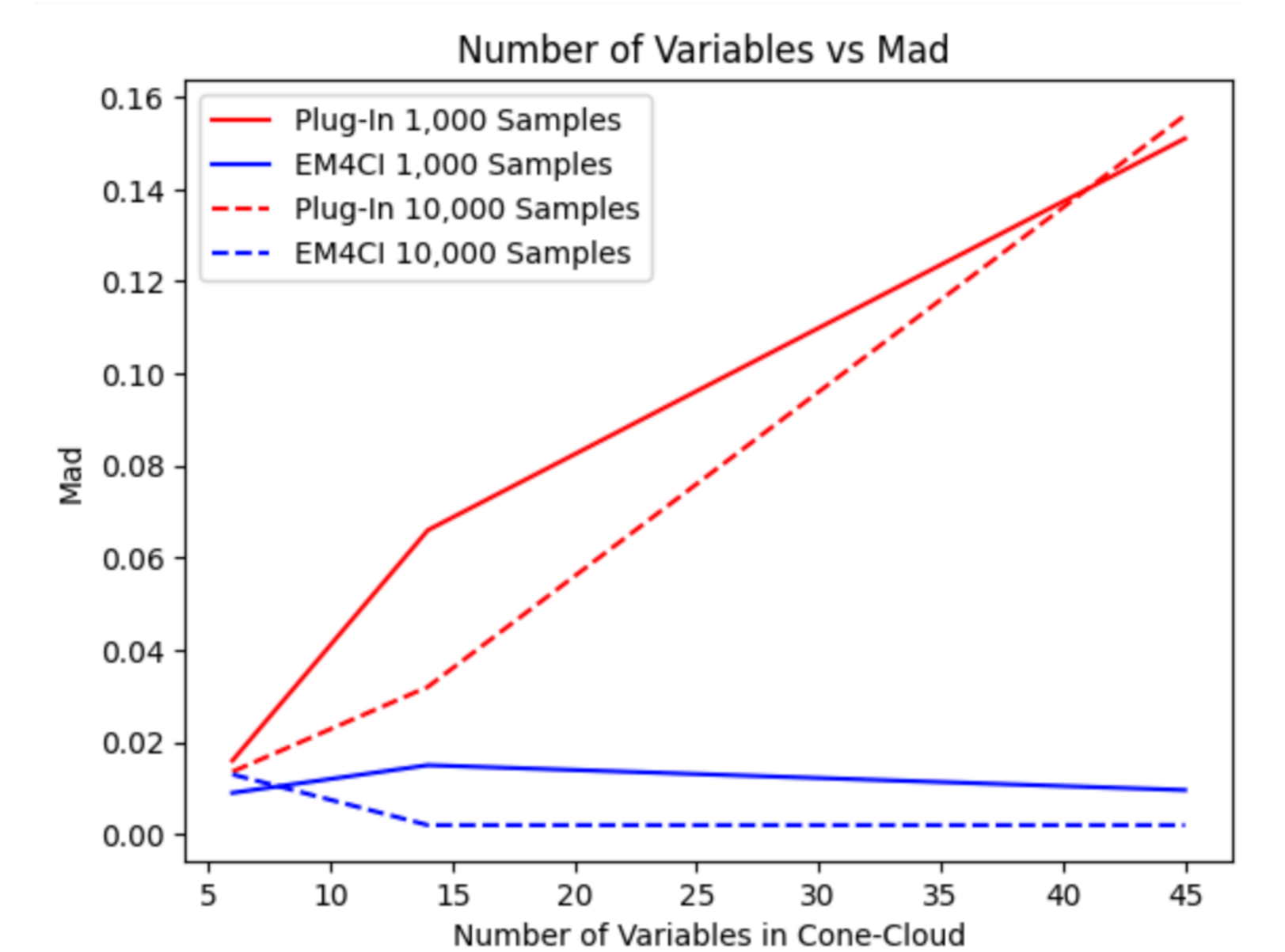
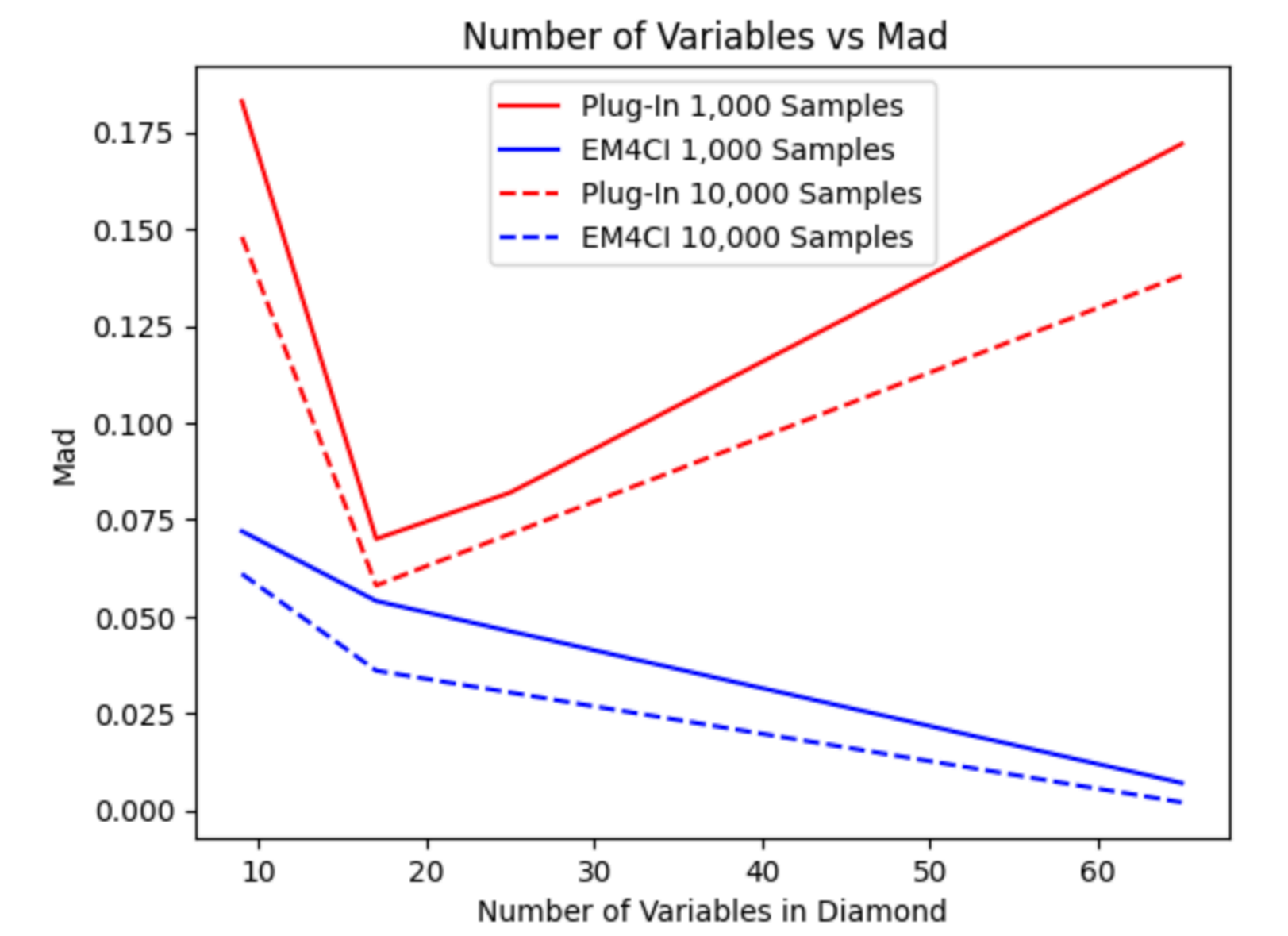
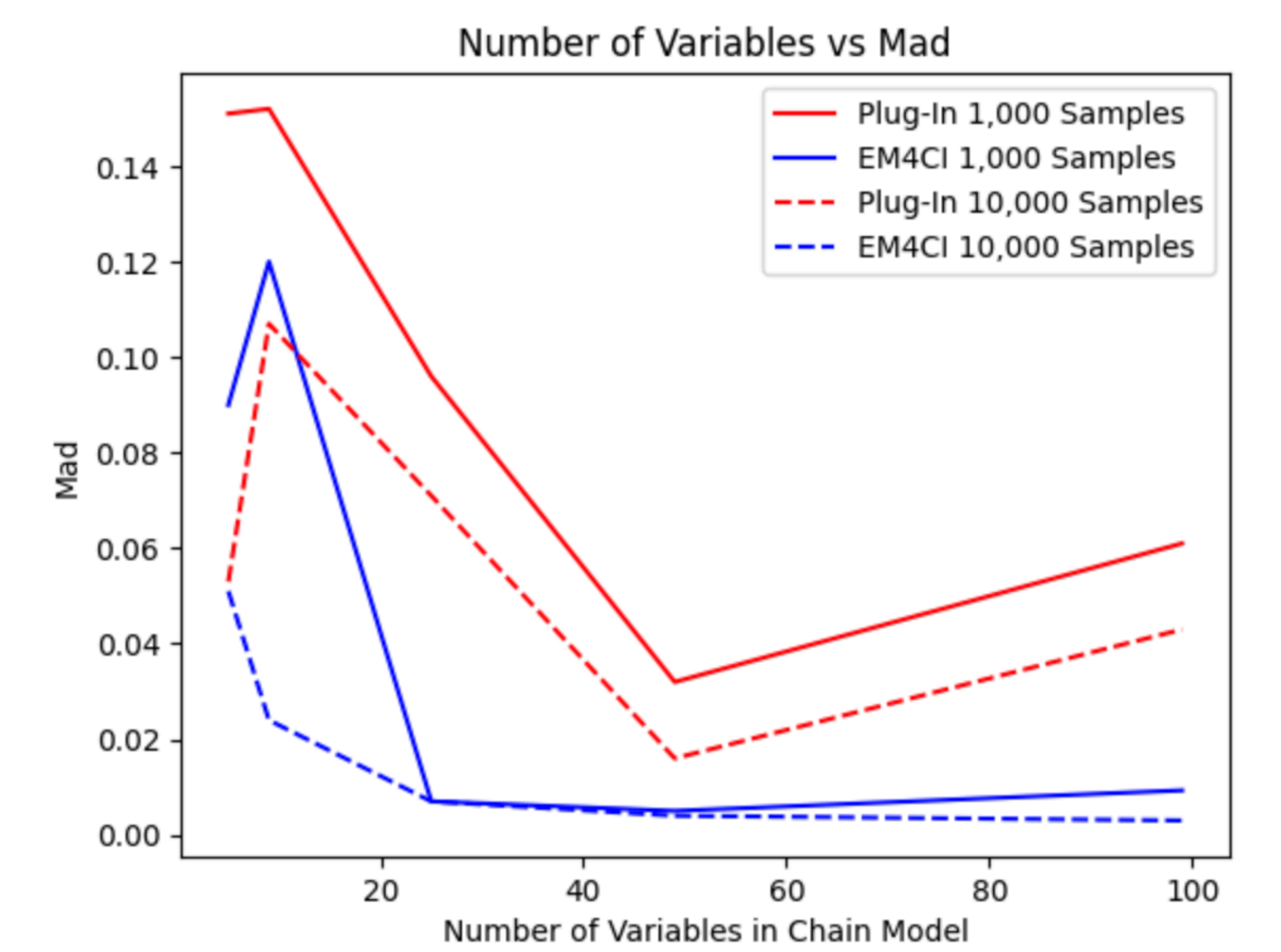
Model	k_{lrm}	100 Samples				1,000 Samples				
		EM4CI mad	EM4CI time(s)	Plug-In mad	Plug-In time(s)	EM4CI mad	EM4CI time(s)	Plug-In mad	Plug-In time(s)	
1	2	0.0037	0.4759	0.0104	1.9	2	0.0032	3.1	0.0025	2.3
2	2	0.1832	1.8643	0.1436	2.3	2	0.0490	8.4	0.0867	2.0
3	2	0.1288	0.9288	0.0569	1.1	2	0.0040	3.6	0.0039	0.7
4	2	0.1819	1.8169	0.1469	2.3	2	0.1438	12.0	0.0704	2.1
5	2	0.4910	1.6539	0.5000	2.0	2	0.0044	17.3	0.0058	2.2
6	2	0.2663	0.3004	0.3930	2.0	2	0.1263	0.5	0.1319	2.1
7	2	0.2520	0.7757	0.2509	1.9	2	0.0891	7.1	0.0238	2.0
8	2	0.1372	0.6348	0.1579	2.0	2	0.2340	4.7	0.1303	1.9

Competing Scheme: WERM [Y. Jung et al., 2020]

Model	k_{lrm}	1,000 Samples				10,000 Samples				
		WERM error	WERM time(s)	EM4CI error	EM4CI time(s)	WERM error	WERM time(s)	EM4CI error	EM4CI time(s)	
1	2	0.0071	18.7	0.0059	8.8	2	0.0031	32.6	0.0046	63.5
8	2	0.1082	25.8	0.1566	7.6	2	0.11	47.7	0.0001	81.4
3'	2	0.027	27.2	0.0004	3.5	2	0.001	44.1	0.0009	53.1

- Learns causal effects by weighted empirical risk minimization.
- State of the art method that focuses on estimating the quantities in the estimand using statistical methods.

Synthetic Network Results

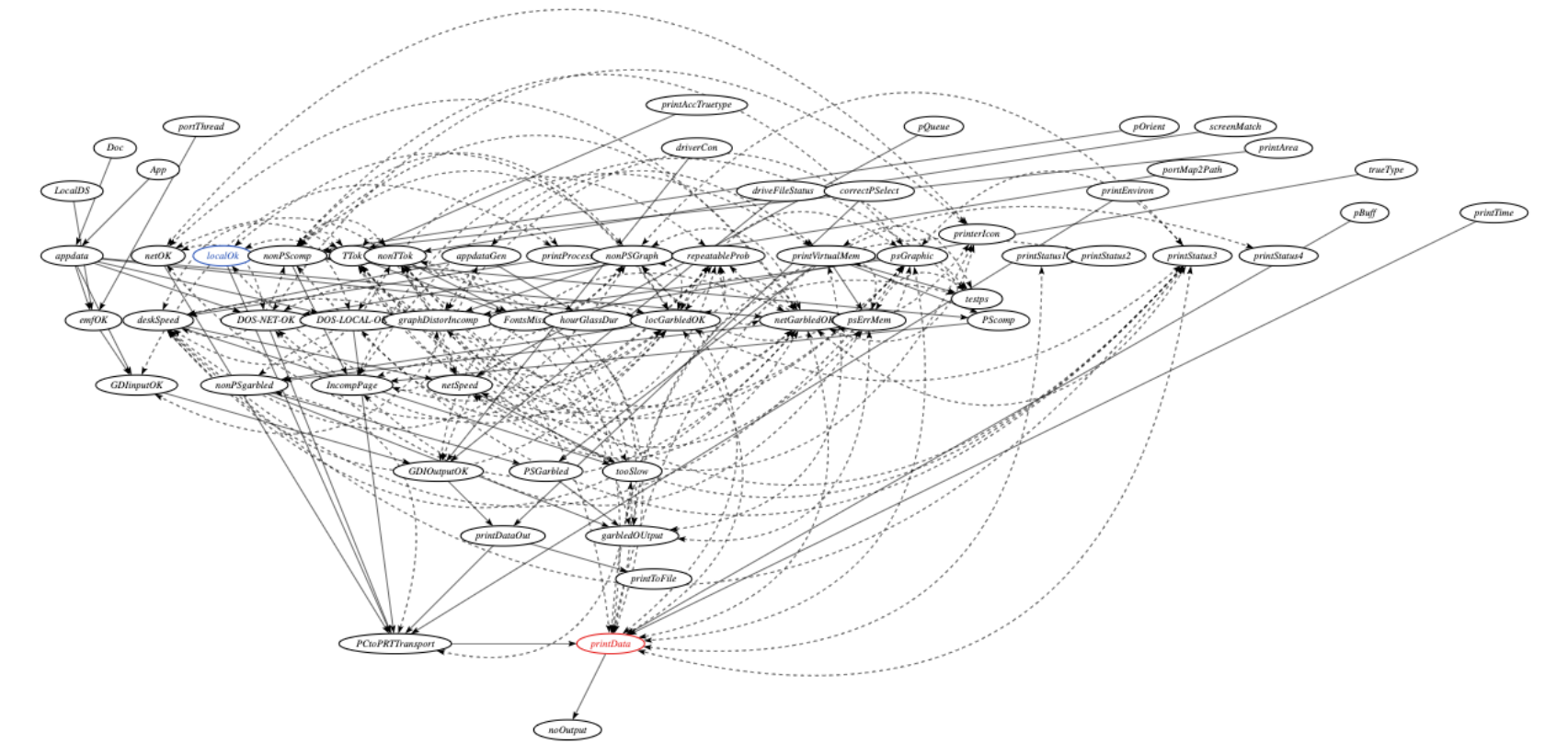


UAI Benchmark Results

Table 1: Plug-In & EM4CI results on the A Network $|V| = 46$; $|U| = 8$; $d = 2$; $k = 2$ treewidth ≈ 16

Query	Plug-In				EM4CI			
	1,000 Samples mad	1,000 Samples time(s)	10,000 Samples mad	10,000 Samples time(s)	1,000 Samples mad	1,000 Samples time(s)	10,000 Samples mad	10,000 Samples time(s)
$P(V_{51} do(V_{10}))$	0.0584	8.0	0.0114	55.7	0.0139	0.0012	0.0083	0.0012
$P(V_{51} do(V_{14}))$	0.0319	8.3	0.0056	51.3	0.0143	0.0047	0.0086	0.0046
$P(V_{51} do(V_{41}))$	0.0255	13.9	0.0092	48.3	0.0147	0.0042	0.0079	0.0041
$P(V_{51} do(V_{45}))$	0.0496	9.8	0.0206	49.1	0.0140	0.0031	0.0082	0.0030

EM4CI Learning time=71(s), $k_{lrm} = 4$ (1,000 Samples) and time=541(s), $k_{lrm} = 4$ (10,000 Samples)



Conclusion

- EM4CI was extremely accurate on all benchmarks we tried.
- Inference on multiple queries was very fast after learning.
- EM4CI is another tool for causal inference, not meant to replace the estimand based approach but used as an alternative when beneficial.