

Causal Inference Despite Limited Global Confounding via Mixture Models

Spencer Gordon ¹

Bijan Mazaheri ¹

Yuval Rabani² Leonard Schulman¹

¹California Institute of Technology

²The Hebrew University of Jerusalem

Problem

Setup

- A Bayesian Network is a directed acyclic graph (DAG) $\mathcal{G} = (\mathbf{V}, \mathbf{E})$.
- ullet A Bayesian Network Distribution on n random variables $\mathbf V$ is Markovian on Bayesian Network $\mathcal G$.
- A k-mixture of such distributions (k-MixBND) is represented using one additional vertex U with CH(U) = V.

Task

Knowns: ullet The marginal probability distribution on ${f V}$:

$\Pr(\mathbf{V}) = \sum \Pr(u) \Pr(\mathbf{V} \mid u)$

Unknowns:

- The probability distribution on U, i.e. Pr(u) for $u \in \{1, \ldots, k\}.$
- The within source probability distribution $\Pr(\mathbf{V} \mid u) = \mathcal{P}_u(\mathbf{V}) \text{ for } u \in \{1, \dots, k\}.$
- $\mathcal{P}_u(\mathbf{V})$ is a Bayesian network distribution, so it suffices to find $\mathcal{P}_u(V \mid \mathbf{pa}(V))$ for all $V \in \mathbf{V}$, and assignments pa(V) to PA(V).

Motivation

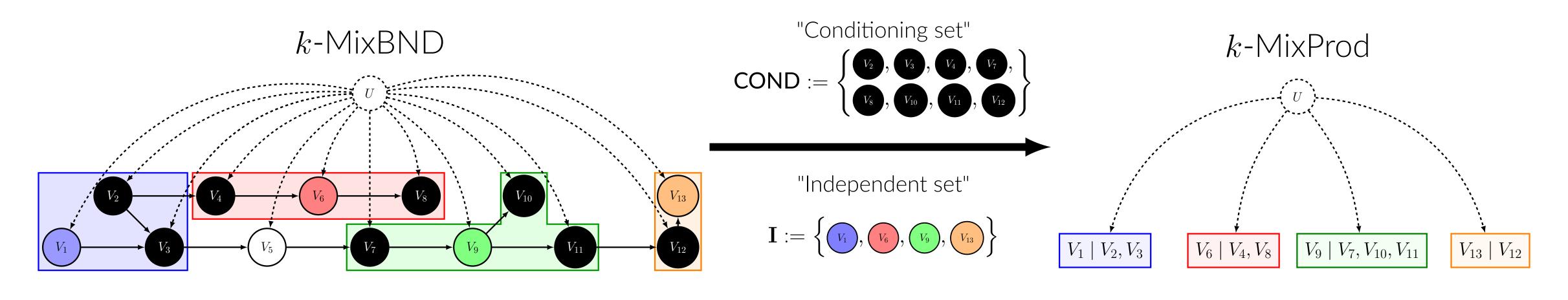
- Setting emerges when combining multiple...
 - -Populations -Environments k < observed support!-Datasets
- Graphically, the causal relationships are considered unidentifiable.
- Interventional probabilities can be calculated using the within source distributions, $Pr(\mathbf{V} \mid u)$.
- Solving the mixture model allows for a degree of deconfounding.

Previous Work

Using 3-independent variabls and larger alphabets (linear in k):

- E. S. Allman, 2009 use algebraic methods to exploit within-source independence.
- Anandkumar, Hsu, and Kakade, 2012 follows a similar strategy using tensor decomposition.
- Wang and Blei, 2019 introduced deconfounders using multiple causes.
- Criticized in Ogburn, Shpitser, and Tchetgen, 2019 and D'Amour, 2019
- Criticism is linked to the difference between learning parameters that generate similar statistics and identification of the true parameters.

Strategy



Reduction to *k***-MixProd**

- The conditional probability distribution I | COND is an instance of k-MixProd on the **independent set I**.
- We will create many such instances (called "runs") and stitch together the results. These runs will need to:
- Cover all of the variables and assignments to their parents. 2. Be alignable with eachother so their results can be synthesized.
- Runtime of k-MixBND: $n2^{\Delta^2}$ executions of k-MixProd where Δ is a bound on the degree of \mathcal{G} .
- Runtime of k-MixProd: $2^{\mathcal{O}(k^2)}n^{\mathcal{O}(k)}$ in Gordon et al., 2021 and $2^{\mathcal{O}(k\log(k))}n^{\mathcal{O}(k)}$ in upcoming work.

Difficulties

- k-MixProd is symmetric to permutations in the labels of U. Solution: Alignment variables
- Conditioning (via post-selection) limits the variables whose information we have access to. How can we ensure we have obtained all of the necessary information?
- Solution: "Good sets of runs" and alignment spanning trees
- The parameters returned from k-MixProd instances are of the form $\mathcal{P}_u(V \mid \mathbf{mb}(V))$ -- we want them of the form $\mathcal{P}_u(V \mid \mathbf{pa}(V))$ (which is standard for Bayesian networks). Solution: Bayesian unzipping

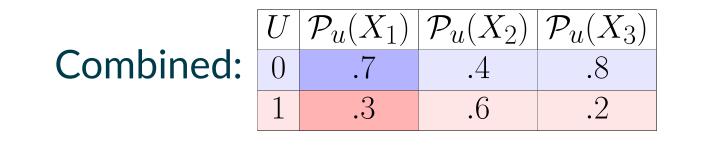
Assumptions

- 1. We have access to a k-MixProd oracle requiring O(k) variables that are independent within each source.
- The observable variables in our BND are binary and discrete.
- Extensions exist for larger alphabets.
- 3. The mixture is supported on $\leq k$ sources.
- 4. The underlying Bayesian DAG is sufficiently sparse.
 - The algorithm works if $n \ge (\Delta + 1)^4 N_{\rm mp}$ • Milder requirements exist for specific graphs.
- 5. The resulting product mixtures are non-degenerate.
- 6. The DAG structure \mathcal{G} is known.
 - Upcoming work on how to do causal discovery to learn \mathcal{G} .

Alignment of Source Labels

$\begin{array}{c|c|c} U^{(a)} & \mathcal{P}_{u^{(a)}}(X_1) & \mathcal{P}_{u^{(a)}}(X_2) \\ \hline 0 & .7 & .4 \\ \hline \end{array}$ $\begin{array}{|c|c|c|c|c|c|} U^{(b)} & \mathcal{P}_{u^{(b)}}(X_1) & \mathcal{P}_{u^{(b)}}(X_3) \\ \hline 0 & .3 & .2 \\ \hline \end{array}$ Run *b*: 0

Alignment Variables

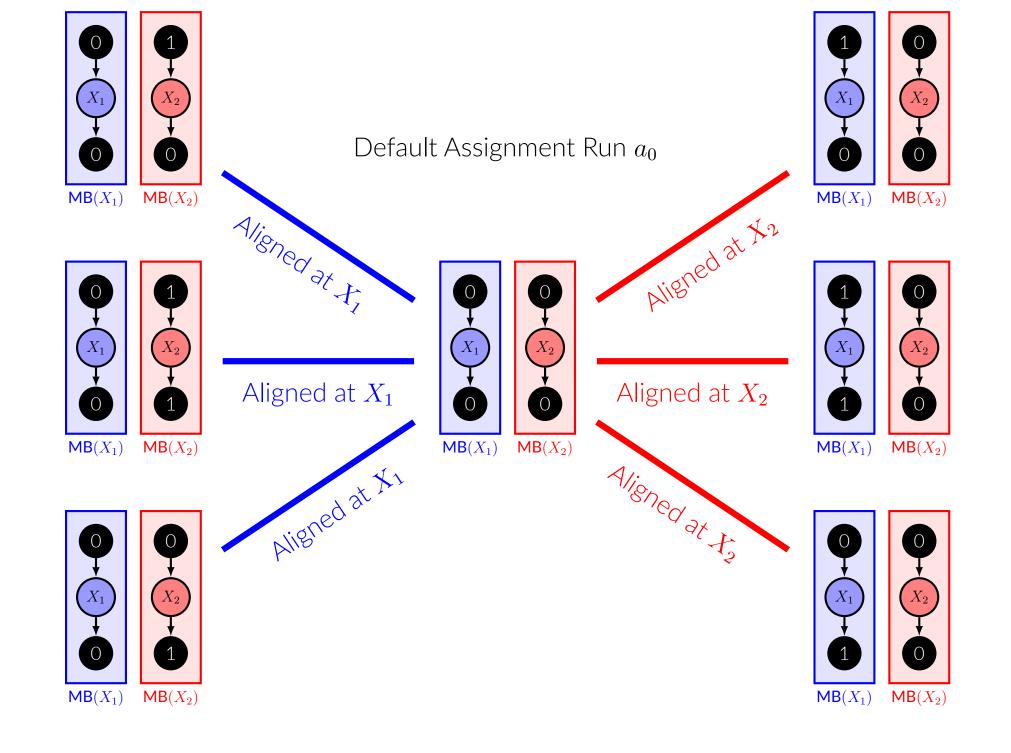


• Runs a and b are aligned at X_1 , allowing $U^{(a)}$ and $U^{(b)}$ to be aligned.

To ensure we have a "good collection of runs," any two runs must be

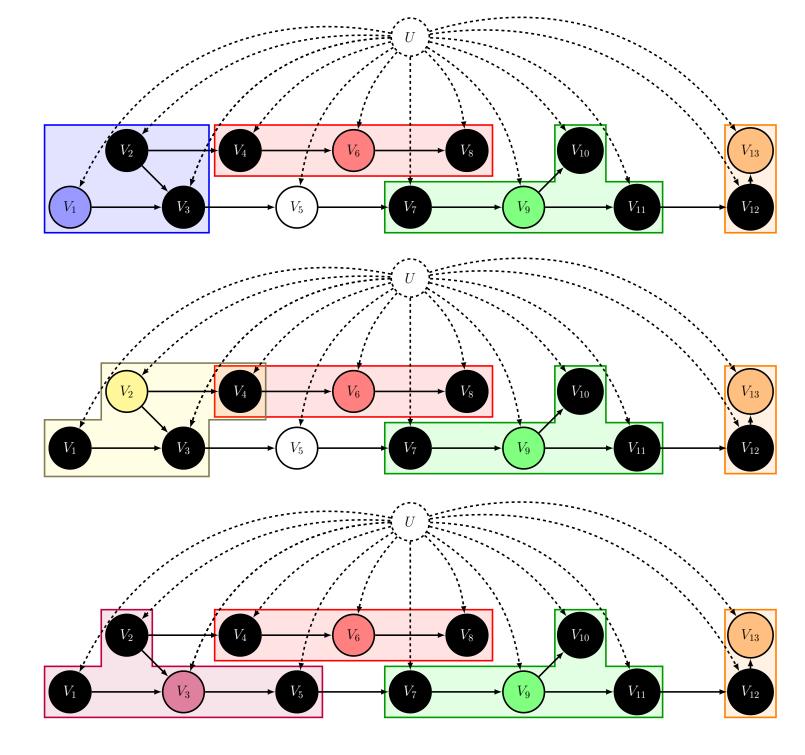
alignable via a chain of alignment variables - i.e. we must have an

Changing the Conditioned Values



Change the assignments to Markov boundaries while leaving at least one assignment the same.

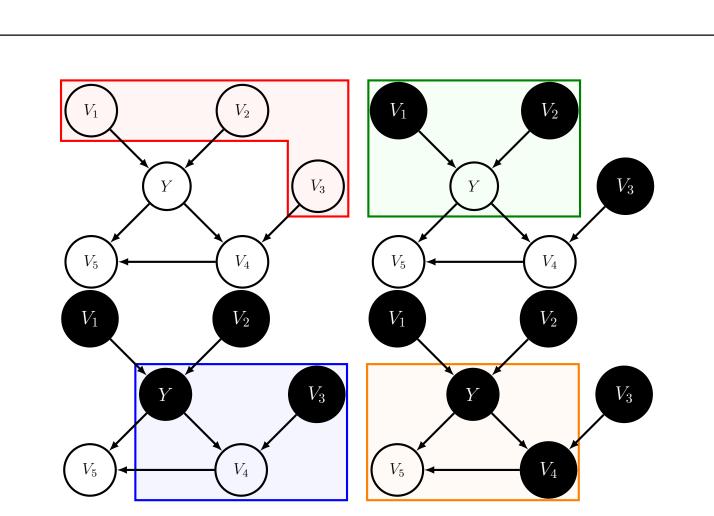
Varying the Independent Set



Disjoint Markov boundaries ensure that a single variable can be swapped without requiring others in the independent set to be conditioned on.

denominator, so it cancels out.

Bayesian Unzipping: $Pr(V \mid MB(V)) \rightarrow Pr(V \mid PA(V))$



alignment spanning tree.

Let y^0 denote ``y=0" and y^1 denote ``y=1".

$$\mathcal{P}_u(y^1\mid \mathsf{mb}(Y)) = \frac{\mathcal{P}_u(y^1,\mathsf{mb}(Y))}{\mathcal{P}_u(y^1,\mathsf{mb}(Y)) + \mathcal{P}_u(y^0,\mathsf{mb}(Y))}$$

We apply the standard factoring,

$$\mathcal{P}_u(y, \mathsf{mb}(Y)) = \mathcal{P}_u(v_1, v_2, v_3) \mathcal{P}_u(y \mid v_1, v_2) \mathcal{P}_u(v_4 \mid y, v_3) \mathcal{P}_u(v_5 \mid y, v_4),$$

to all three terms.

$$\mathcal{P}_{u}(y^{1} \mid \mathbf{mb}(Y)) = \frac{\mathcal{P}_{u}(y^{1} \mid v_{1}, v_{2}) \mathcal{P}_{u}(v_{4} \mid y^{1}, v_{3}) \mathcal{P}_{u}(v_{5} \mid y^{1}, v_{4})}{\mathcal{P}_{u}(y^{1} \mid v_{1}, v_{2}) \mathcal{P}_{u}(v_{4} \mid y^{1}, v_{3}) \mathcal{P}_{u}(v_{5} \mid y^{1}, v_{4})) + \mathcal{P}_{u}(y^{0} \mid v_{1}, v_{2}) \mathcal{P}_{u}(v_{4} \mid y^{0}, v_{3}) \mathcal{P}_{u}(v_{5} \mid y^{0}, v_{4})}$$

- $\mathcal{P}_u(v_1, v_2, v_3)$ appears in both numerator and
- If we traverse in reverse topological order, then $\mathcal{P}_u(v_4 \mid y, v_3)$ and $\mathcal{P}_u(v_5 \mid y, v_4)$ terms are previously calculated for both $y \in \{y^0, y^1\}$.
- $\mathcal{P}_u(y^0 \mid v_1, v_2) + \mathcal{P}_u(y^1 \mid v_1, v_2) = 1$, so we can solve for green terms.
- Iterating this process incurs stability costs proportional to the depth of the graph.
- This can be avoided by not conditioning on the children of the deepest variables in the independent set.

References

(2012). ``A Method of Moments for Mixture Models and Hidden Markov Models". In: *Proc. 25th Ann. Conf. on Learn-* 114.528, pp. 1597–1601. ceedings, pp. 33.1-33.34. URL: http: // proceedings . mlr . press / v23 / anandkumar12/anandkumar12.pdf.

Anandkumar, A., D. J. Hsu, and S. M. Kakade D'Amour, Alexander (2019). ``Comment: Reflections on the deconfounder". In: Journal of the American Statistical Association Gordon, S. L. et al. (2021). ``Source Identi-Ogburn, Elizabeth L, Ilya Shpitser, and Eric

ing Theory - COLT. Vol. 23. JMLR Pro- E. S. Allman C. Matias, J. A. Rhodes (2009). 'Identifiability of parameters in latent structure models with many observed vari-

ables". In: Ann. Statist. 37.6A, pp. 3099-3132. DOI: 10.1214/09-A0S689.

fication for Mixtures of Product Distributions". In: Proc. 34th Ann. Conf. on Learning Theory - COLT. Vol. 134. Proc. Machine Learning Research. PMLR, pp. 2193-

2216. URL: http://proceedings.mlr.Wang, Y. and D. M. Blei (2019). ``The press/v134/gordon21a.html.

J Tchetgen Tchetgen (2019). ``Comment on "blessings of multiple causes"". In: Journal of the American Statistical Association 114.528, pp. 1611-1615.

Blessings of Multiple Causes". In: Journal of the American Statistical Association 114.528, pp. 1574-1596. DOI: 10.1080/ 01621459.2019.1686987.