

Learning Discrete Directed Acyclic Graphs via Backpropagation

Andrew J. Wren[∇] Pasquale Minervini^{∇Ω}
Luca Franceschi[∏] Valentina Zantedeschi^{∗∇∩}

andrew.wren@ntlworld.com

arXiv:2210.15353

Andrew J. Wren is very grateful to his co-authors who supervised him in the University College London MSc Machine Learning thesis project from which this paper derives.

[∇] University College London ^Ω University of Edinburgh [∩] Inria London

[∏] Amazon Web Services (work done while at UCL, and prior to joining Amazon)

[∗] ServiceNow

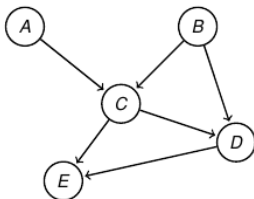


servicenow[®]



- Problem: learning a Bayesian network's DAG from data.
- Question: can a probabilistic discrete backpropagation approach be used?
- Backprop methods.
- Architecture: DAG-DB.
- Experiments and results.
- Conclusion and future work.

- Recall the problem of structure identification: learning a Bayesian network's directed acyclic graph (DAG) from data it has generated.



$$P(A, B, C, D, E) = P(A)P(B)P(C|A, B)P(D|B, C)P(E|C, D)$$

Figure. Example of Bayesian network, with $d = 5$ nodes.^[1]

- Given data $\mathbf{X} = (\mathbf{x}_{n,A}, \dots, \mathbf{x}_{n,E})_{n=1}^N \in \mathbb{R}^{N \times d}$ from a hidden d -node DAG D , make prediction D_{pred} .

[1]Adapted from Sahani, UCL (2021).

Problem: learning a DAG from data (2)

- Existing approaches to the problem can mainly be divided between combinatoric, continuous and Bayesian.
- Differences include representation of DAG edges in core calculations, and whether gradient descent (GD) is used.

Approach	DAG edges	GD?	Examples
Combinatoric	Binary	No	PC-Stable ^[1] , FGES ^[2]
Continuous	Float	Yes	NOTEARS ^[3] , GOLEM ^[4]
Bayesian	Binary	Yes	DiBS ^[5] , DAG-GFlowNet ^[6]
Probabilistic	Binary	Yes	<u>DAG-DB</u> <i>presented today</i>

[1]Colombo & Maathuis. *Order-independent constraint-based causal structure learning*. (2014)

[2]Ramsey, Glymour, Sanchez-Romero & Glymour. *A million variables and more...* . (2017)

[3]Zheng, Aragam, Ravikumar & Xing. *DAGs with NO TEARS....* (2018)

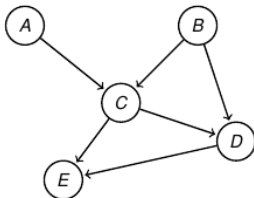
[4]Ng, Ghassami & Zhang. *On the Role of Sparsity and DAG Constraints for Learning Linear DAGs*. (2020)

[5]Lorch, Rothfuss, Schölkopf & Krause. *DiBS: Differentiable Bayesian Structure Learning*. (2021)

[6]Dele, Góis, Emezue, Rankawat, Lacoste-Julien, Bauer & Bengio. *Bayesian Structure Learning with Generative Flow Networks*. (2022)

Underlying approach:

- For data point \mathbf{x} , predict value at a 'child' node from its 'parents'.



$$P(A, B, C, D, E) = P(A)P(B)P(C|A, B)P(D|B, C)P(E|C, D)$$

Figure. Example of Bayesian network.^[1]

- E.g, predict \mathbf{x}_D from \mathbf{x}_B and \mathbf{x}_C .

[1]Adapted from Sahani, UCL (2021).

Framework:

- **Directed graph.** With $\mathbf{X} \in \mathbb{R}^{N \times d}$, looking for graph on d nodes. Let \mathbf{Z} be $d \times d$ *binary* adjacency matrix for a directed graph.
- **Distribution.** A *float* matrix Θ (same shape as \mathbf{Z}) parametrizes an ExpFam distribution,

$$p(\mathbf{Z}; \Theta) = \exp(\langle \mathbf{Z}, \Theta \rangle - A(\Theta)).$$

- For **sampling**, use ‘Perturb-&-MAP’^[1] to get S samples. For each sample $\mathbf{Z}^{(s)}$, take the most probable value of \mathbf{Z} given parameter Θ plus noise:

$$\mathbf{Z}^{(s)} = \text{MAP}(\Theta + \epsilon^{(s)}) \text{ with } \epsilon^{(s)} \sim \text{Logistic}(0, 1), \quad s = 1, \dots, S.$$

[1]Papandreou & Yuille. *Perturb-and-MAP random fields: Using discrete optimization to learn and sample from energy models.* (2011)

To learn Θ , need a way to **backpropagate** from \mathbf{Z} to Θ .
Backprop methods which keep \mathbf{Z} binary include:

- **Straight-Through Estimation (STE)**.^[1] Very simple.
- **Score-Function Estimation (SFE)**. Aka REINFORCE.
- **Blackbox Estimation (BBE)**. Not probabilistic.
- **Implicit Maximum Likelihood Estimation (I-MLE)**.^[2]
Taking advantage $p(\mathbf{Z}; \Theta) \in \text{ExpFam}$.

^[1]Hinton, Srivastava & Swersky. *Neural networks for machine learning*. (2012); Bengio, Léonard & Courville. *Estimating or Propagating Gradients Through Stochastic Neurons for Conditional Computation*. (2013)

^[2]Niepert, Minervini & Franceschi. *Implicit MLE: Backpropagating Through Discrete Exponential Family Distributions*. (2021)

- For **backward pass**, with average loss (plus regularizers) L , Straight-Through Estimation (STE), approximates gradient of L w.r.t. Θ as

$$\frac{\partial L}{\partial \Theta} \propto \sum_{s=1}^S \frac{\partial L}{\partial \mathbf{z}^{(s)}}.$$

- Or, for **backward pass** with Implicit Maximum Likelihood (I-MLE), perturb Θ using each sample $\mathbf{z}^{(s)}$,

$$\tilde{\Theta}^{(s)} = \Theta - \lambda \frac{\partial L}{\partial \mathbf{z}^{(s)}}.$$

- Use noise values from sampling $\epsilon^{(s)}$ to set

$$\tilde{\mathbf{z}}^{(s)} = \text{MAP}(\tilde{\Theta}^{(s)} + \epsilon^{(s)}).$$

- Approximate gradient of L w.r.t. Θ as

$$\frac{\partial L}{\partial \Theta} \approx \frac{1}{\lambda S} \sum_{s=1}^S \left[\mathbf{z}^{(s)} - \tilde{\mathbf{z}}^{(s)} \right].$$

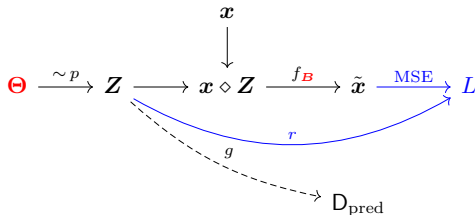
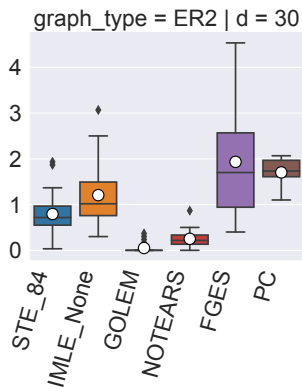


Figure. DAG-DB, including **learnable parameters** and **loss calculation**.

- $x \diamond Z$ ensures only ‘parents’ influence predictions for ‘children’.
- f_B is a very simple neural net: linear no bias.
- r is regularizer: NOTEARS-like regularizer pushing Z towards being a DAG, plus sparsity regularizer.
- Depending on compute and size of problem, Z is transformed into a DAG D_{pred} by a maximum DAG algorithm g in training or only at evaluation.



ER2: linear Gaussian additive noise model on random (Erdős-Rényi) graphs with d nodes, and expected number of edges is $2d$. In this example, $d = 30$ nodes.

Results on synthetic data showing a normalised class **structural Hamming distance** $nSHD_c$ from the true DAG.

Shows two DAG-DB methods STE_84 and IMLE_None, compared with continuous (GOLEM, NOTEARS) and combinatoric (FGES, PC) methods.

Model		SHD _c	prec _c
IMLE_None	mean	12.7	0.869
	median	13	1.000
GOLEM		11	1.000
NOTEARS		11	0.467
FGES		11	0.750
PC		11	0.750

Tests on the Sachs cellular biochemistry dataset.¹ Best metric scores are in bold. SHD_c is un-normalised; prec_c is a **precision** metric.

¹ Sachs, Perez, Pe'er, Lauffenburger & Nolan. *Causal protein-signaling networks derived from multiparameter single-cell data.* (2005)

Take home

- Using discrete probabilistic backpropagation, DAG-DB performs competitively at structure identification, often better than combinatoric methods, though not as well as continuous approaches.
- May be easier than many other DAG identification methods to integrate with other neural nets.
- Extend to more types of Bayesian network: non-Gaussian, non-linear, discrete-valued, with interventions.
- What if no 'gold' DAG, e.g. we generate a latent DAG optimised to achieve some goal?

Thank you!