

# Exploration of Losses in NOTEARS

Mengqi Liu

Aug 30, 2023

## 1 Notations

NOTEARS<sup>[1]</sup> solves the problem of structure learning for directed acyclic graph(DAG) under the assumption that data is generated by generalized linear SEM. Consider a vector of random variables of our interest  $X = (X_1, X_2, \dots, X_d)$ . Let  $\mathbf{X} \in \mathbb{R}^{n \times d}$  be the data matrix, where  $(i, j)$ -th entry  $X_{ij}$  represents the  $i$ -th i.i.d observation of the  $j$ -th variable  $X_j$ .

Under the assumption that the data fits the generalized linear SEM, we could model  $X$  through a weighted adjacency matrix  $W = (w_1, w_2, \dots, w_d) \in \mathbb{R}^{d \times d}$  by  $\mathbb{E}(X_j) = f(w_j^\top X)$  for  $j = 1, 2, \dots, d$ .

Note  $A(W) \in \{0, 1\}^{d \times d}$  as the adjacency matrix with weights in  $W$  removed. Specifically, let  $[A(W)]_{ij}$  be 0 if  $W_{ij} = 0$  and 1 otherwise. Let  $G(W)$  denote the graph corresponding to adjacency matrix  $A(W)$ . Also let  $\mathbb{D}$  denote the discrete space of DAGs and  $G = (V, E)$  denote a graph with vertices  $V$  and edges  $E$ . Define the number of edges  $s_0 = |E|$ . Our causal discovery task is to learn a DAG  $G \in \mathbb{D}$  given  $\mathbf{X}$ . Under settings of Bayesian networks, we could model  $X$  in this way:  $\mathbb{E}(X_j | \text{pa}(X_j)) = f(w_j^\top X)$  for  $j = 1, 2, \dots, d$ . And notice  $W_{ij} = 0$  for any  $i \in \{1, 2, \dots, d\} \setminus \text{pa}(X_j)$ .

Especially, we focus on linear SEM first, where  $X_j = w_j^\top X + \epsilon_j$  for  $j = 1, 2, \dots, d$ . Here we do not restrict  $\epsilon_j$ 's to Gaussian noises.

## 2 Review on losses in NOTEARS

In the paper, Xun Zheng et al. use L2 loss to score DAGs under linear SEM. The benefits are as follows: (1) The minimizer of the L2 loss has been shown to recover a true DAG with high probability on finite-samples and in high-dimensions. (2) The estimator through the L2 loss is consistent for both Gaussian SEM and non-Gaussian SEM. (3) No faithfulness assumption needed in the setting of linear SEM and L2 loss. The L2 loss under linear SEM is

$$\ell(W; \mathbf{X}) = \frac{1}{2n} \|\mathbf{X} - \mathbf{X}W\|_F^2. \quad (1)$$

Due to the interest in learning sparse DAG, authors introduced a  $\ell_1$ -regularization  $\|W\|_1 = \|\text{vec}(W)\|_1$  in the regularized score function:

$$F(W) = \ell(W; \mathbf{X}) + \lambda \|W\|_1 = \frac{1}{2n} \|\mathbf{X} - \mathbf{X}W\|_F^2 + \lambda \|W\|_1, \quad \lambda \geq 0. \quad (2)$$

They put forward a smooth function  $h(W)$  to characterize the topological property of  $G(W)$  and have proven that  $h(W) = 0$  if and only if  $W$  is acyclic (i.e.  $G(W) \in \mathbb{D}$ ).

$$h(W) = \text{tr}(e^{W \circ W}) - d = 0 \quad (3)$$

Then the causal discovery question becomes a equality-constrained program(ECP) as

$$\begin{aligned} \min_{W \in \mathbb{R}^{d \times d}} F(W) \\ \text{subject to } h(W) = 0. \end{aligned} \quad (4)$$

In order to solve this ECP, they take advantage of augmented Lagrangian method to solve the original problem by a approximate program:

$$\begin{aligned} \min_{W \in \mathbb{R}^{d \times d}} O(W, \rho) \\ \text{subject to } h(W) = 0, \end{aligned} \quad (5)$$

where  $\rho > 0$  and

$$\begin{aligned} O(W, \rho) &= F(W) + \frac{\rho}{2} |h(W)|^2 \\ &= \frac{1}{2n} \|\mathbf{X} - \mathbf{X}W\|_F^2 + \lambda \|W\|_1 + \frac{\rho}{2} [\text{tr}(e^{W \circ W}) - d]^2. \end{aligned} \quad (6)$$

And the dual problem is

$$\max_{\alpha \in \mathbb{R}} \min_{W \in \mathbb{R}^{d \times d}} D(W, \rho, \alpha) \quad (7)$$

where

$$\begin{aligned} D(W, \rho, \alpha) &= O(W, \rho) + \alpha h(W) \\ &= \frac{1}{2n} \|\mathbf{X} - \mathbf{X}W\|_F^2 + \lambda \|W\|_1 + \frac{\rho}{2} [\text{tr}(e^{W \circ W}) - d]^2 + \alpha [\text{tr}(e^{W \circ W}) - d]. \end{aligned} \quad (8)$$

And the problem was numerically solved by dual ascent method (iteratively update  $\alpha$  and  $\rho$ ) and L-BFGS ( $\lambda = 0$ ) / PQN ( $\lambda > 0$ , find approximation of descent direction for non-smooth function).

### 3 Adversarial relationship between $\ell(W; \mathbf{X})$ and $h(W)$

Notice two facts:

(1)  $W = I_d$  is always a solution of  $\ell(W; \mathbf{X}) = 0$ . But acyclicity requires that diagonal of  $W$  should be 0.

(2) If there's linear relationship between  $\tilde{x}_i$  and  $\tilde{x}_j$ ,  $W_{ij}$  and  $W_{ji}$  should play the same role in  $\ell(W; \mathbf{X})$  while acyclicity needs at least one of the two to be 0.

Therefore intuitively minimizing  $L^\rho(W, \alpha)$  will lead to a trade-off between  $\ell(W; \mathbf{X})$  and  $h(W)$ .

## 4 Theoretical and empirical global minimum for L2 loss

### 4.1 Propositions

**Lemma 4.1** Assume  $W$  is lower triangular with zero diagonal w.o.l.g. (i.e. variables have been arranged according to their topological order),

$$\min_{G(W) \in DAG} \ell(W; \mathbf{X}) = \min_{G(W) \in DAG} \frac{1}{2n} \sum_{j=\max\{d-n+1, 1\}}^d \|\mathbf{X}_j - \mathbf{X}_j w_j\|_2^2.$$

**Proof:** When  $d \leq n$ ,  $\max\{d-n+1, 1\} = 1$  and the statement holds naturally. When  $d > n$ , for any  $j \in \{1, 2, \dots, d-n\}$  (here we choose  $j > 1$  w.l.o.g.), consider the inhomogeneous (probably) linear system with  $W$  as unknown:

$$\begin{pmatrix} X_{11} & X_{12} & \cdots & X_{1,j-1} & X_{1,j+1} & \cdots & X_{1d} \\ X_{21} & X_{22} & \cdots & X_{2,j-1} & X_{2,j+1} & \cdots & X_{2d} \\ \vdots & \vdots & & \vdots & \vdots & & \vdots \\ X_{n1} & X_{n2} & \cdots & X_{n,j-1} & X_{n,j+1} & \cdots & X_{nd} \end{pmatrix} \begin{pmatrix} 0 \\ \vdots \\ 0 \\ W_{j+1,j} \\ \vdots \\ W_{dj} \end{pmatrix} = \begin{pmatrix} X_{1j} \\ X_{2j} \\ \vdots \\ X_{nj} \end{pmatrix}.$$

The linear system can be truncated to

$$\begin{pmatrix} X_{11} & X_{12} & \cdots & X_{1,d-j-1} \\ X_{21} & X_{22} & \cdots & X_{2,d-j-1} \\ \vdots & \vdots & & \vdots \\ X_{n1} & X_{n2} & \cdots & X_{n,d-j-1} \end{pmatrix} \begin{pmatrix} W_{j+1,d} \\ W_{j+2,d} \\ \vdots \\ W_{d,j} \end{pmatrix} = \begin{pmatrix} X_{1j} \\ X_{2j} \\ \vdots \\ X_{nj} \end{pmatrix},$$

where  $d - (j+1) + 1 \geq n$  and  $W_{j+1,j}, W_{j+2,j}, \dots, W_{d,j}$  can take any value in  $\mathbb{R}$ .

Because rows of  $\mathbf{X}$  stand for  $n$  i.i.d observations, we know that  $\text{rank}(\mathbf{X}) = n$  and the corresponding augmented matrix also should have rank of  $n$ . According to Rouché–Capelli theorem, as the rank of its coefficient matrix  $\mathbf{X}$  is equal to the rank of its augmented matrix, this system has at least one solution. Therefore  $\min_{G(W) \in \text{DAG}} \frac{1}{2n} \|\mathbf{X}_j - \mathbf{X}_j w_j\|_2^2 = 0$  for all  $1 \leq j \leq d-n$ .

As a result,

$$\begin{aligned} \min_{G(W) \in \text{DAG}} \ell(W; \mathbf{X}) &= \min_{G(W) \in \text{DAG}} \frac{1}{2n} \|\mathbf{X} - \mathbf{X}W\|_F^2 \\ &= \min_{G(W) \in \text{DAG}} \frac{1}{2n} \sum_{j=1}^d \|\mathbf{X}_j - \mathbf{X}_j w_j\|_2^2 \\ &= \frac{1}{2n} \sum_{j=1}^d \min_{w_j} \|\mathbf{X}_j - \mathbf{X}_j w_j\|_2^2 \\ &= \min_{G(W) \in \text{DAG}} \frac{1}{2n} \sum_{j=\max\{d-n+1, 1\}}^d \|\mathbf{X}_j - \mathbf{X}_j w_j\|_2^2. \end{aligned}$$

□

**Lemma 4.2** Let  $\pi$  denote a permutation and  $P$  is the corresponding permutation matrix. Theoretically,

$$\min_{G(W) \in \text{DAG}} \mathbb{E} \|\mathbf{X} - W^\top \mathbf{X}\|_2^2 \geq \min_P \min_{G(W) \in \text{DAG}} \mathbb{E} \|P\mathbf{X} - W^\top P\mathbf{X}\|_2^2.$$

Empirically,

$$\begin{aligned} \min_{G(W) \in \text{DAG}} \ell(W; \mathbf{X}) &= \min_{G(W) \in \text{DAG}} \frac{1}{2n} \sum_{j=1}^d \|\mathbf{X}_j - \mathbf{X}_j w_j\|_2^2 \\ &\geq \min_P \min_{G(W) \in \text{DAG}} \frac{1}{2n} \sum_{j=1}^d \|\mathbf{X}_j P^\top - \mathbf{X} P^\top w_j\|_2^2 \\ &= \min_P \min_{G(W) \in \text{DAG}} \ell(W; \mathbf{X} P^\top). \end{aligned}$$

**Remark.** Notice that  $\arg \min_{G(W) \in \text{DAG}} \mathbb{E} \|PX - W^\top PX\|_2^2 \neq \arg \min_{G(W) \in \text{DAG}} \mathbb{E} \|X - W^\top X\|_2^2$  in general. Consider the linear regression  $X = W^\top X + \epsilon$ , where  $X \in \mathbb{R}^d$ ,  $W \in \mathbb{R}^{d \times d}$ ,  $\epsilon \in \mathbb{R}^d$ ,  $\text{cov}(X) = \Sigma$  and  $\text{cov}(\epsilon) = D$ . Assume  $W$  is lower triangular, we could conduct Cholesky decomposition as

$$\Sigma = (I - W^\top)^{-1} D (I - W^\top)^{-1} = [(I - W^\top)^{-1} D^{\frac{1}{2}}] [D^{\frac{1}{2}} (I - W^\top)^{-1}] = L(X) L(X)^\top.$$

After permutation, the decomposition becomes

$$P \Sigma P^\top = P L(X) L(X)^\top P^\top = L(PX) L(PX)^\top,$$

where  $L(PX)$  corresponds to new  $\tilde{W}$  with respect to  $PX = \tilde{W}^\top PX + P\epsilon$ . Generally,  $\tilde{W} = I - \{L(PX)[PDP^\top]^{-\frac{1}{2}}\}^{-1}$  is different from original  $W = I - \{L(X)D^{-\frac{1}{2}}\}^{-1}$ .

## 4.2 Empirical results

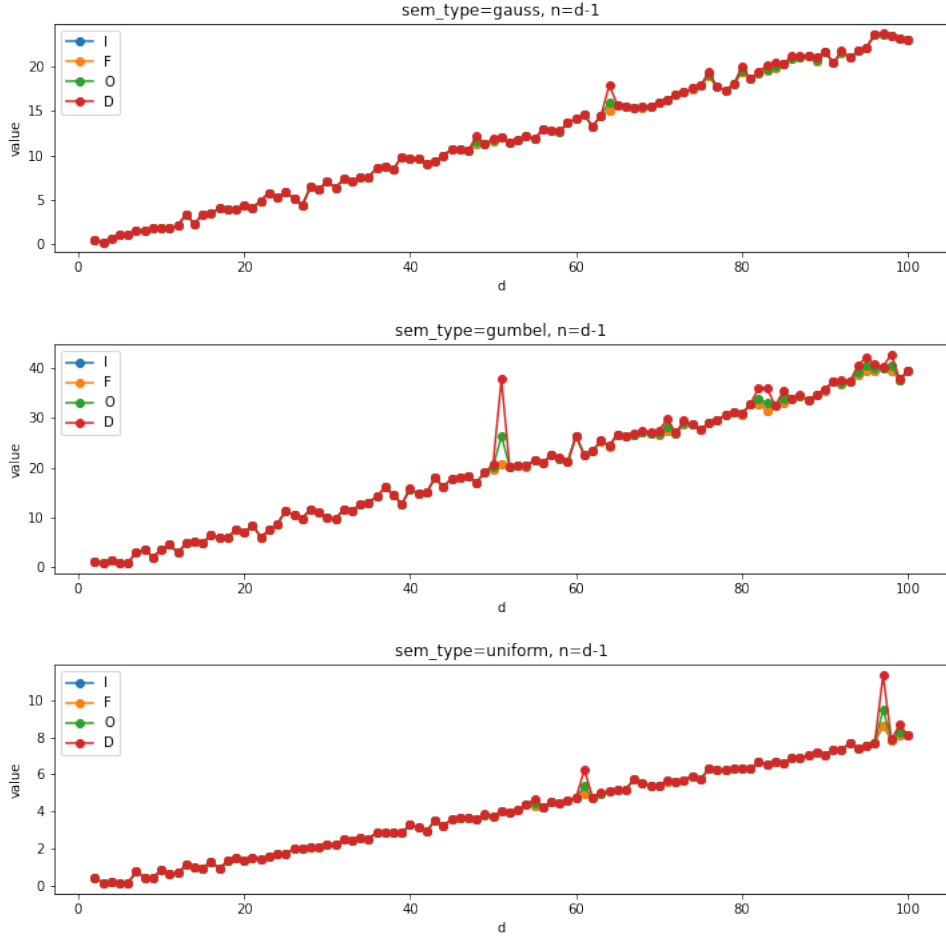


Figure 1:  $d - n = 1$

To check the performance of empirical L2 loss, we only train DAGs with data generated with noises of continuous distributions (Gaussian, Uniform, and Gumbel in total) and set  $\lambda = 0$ .

First, we let  $d - n = 1$ ,  $s_0 = d - 1$  and train with  $d \in \{2, 3, \dots, 100\}$  respectively. The empirical losses are shown in Figure 1. From the results, we could find an increasing trend in L2 loss with increasing  $d$ (or  $n$ ). This is reasonable because larger  $d$  will introduce more noise and the right side of Lemma 4.1 will increase accordingly.

Second, zooming in and observing losses when  $d \leq 20$  in Figure 2, where different lines refer to repeated experiments with different random seeds. random fluctuations when  $d$  is relatively small are not catastrophic to subvert the trend for now.

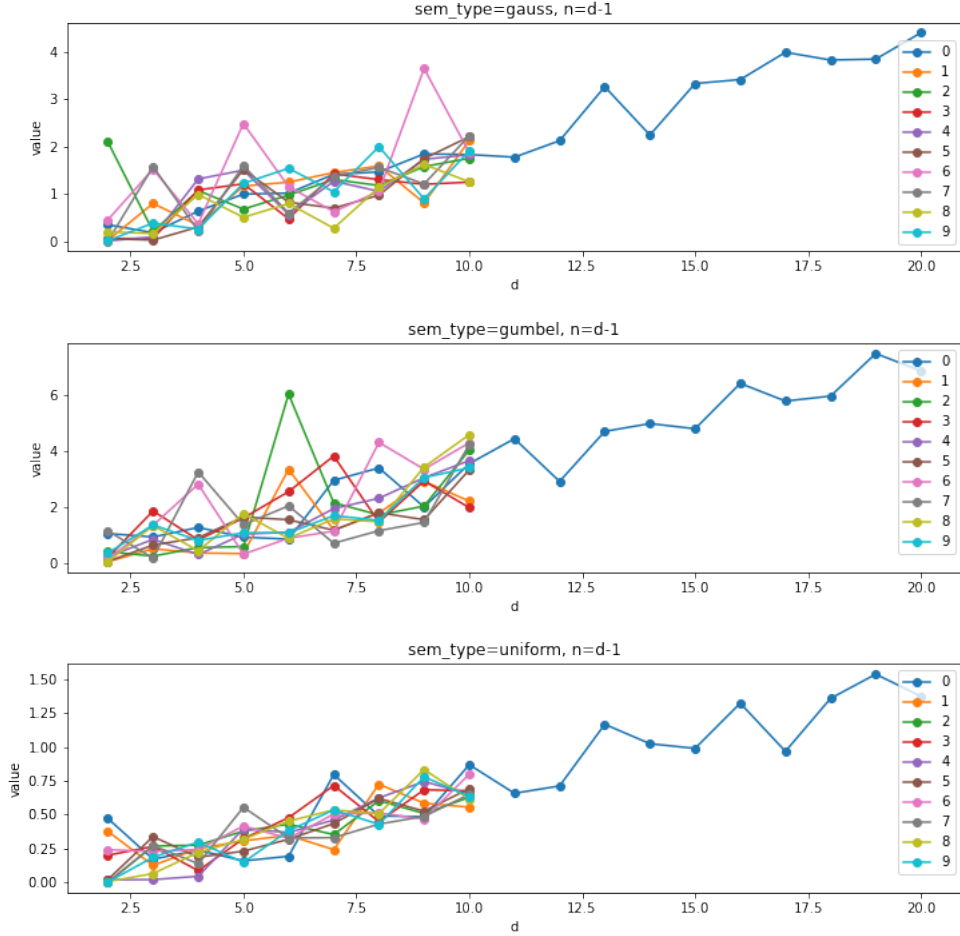


Figure 2:  $d \leq 20$

Third, we try to permute the columns of data matrix  $\mathbf{X}$  to construct a lower bound of empirical L2 loss and to verify if NOTEARS has attained the empirical global minimum. The losses optimized by L-BFGS are showed in Figure 3. We could find that all empirical minimized L2 losses are almost the same (although they're not exactly the same after zooming in). But it is worth noting that the gap between permutations become larger when  $d$  increases. Similarly, we could find the learned losses are very close with permutations in Figure 4. It is worth noting that for SGD, I set the learning rate to be  $\frac{1}{\rho}$  avoid infinite loss that occurred in the previous experiments. Also I only test for Gaussian distributions and fewer  $d$ 's in SGD due to the time-consuming learning process.

This can be showed as an initial trial to verify that NOTEARS has performed good optimization with respect to L2 loss through L-BFGS as it originally put forward.

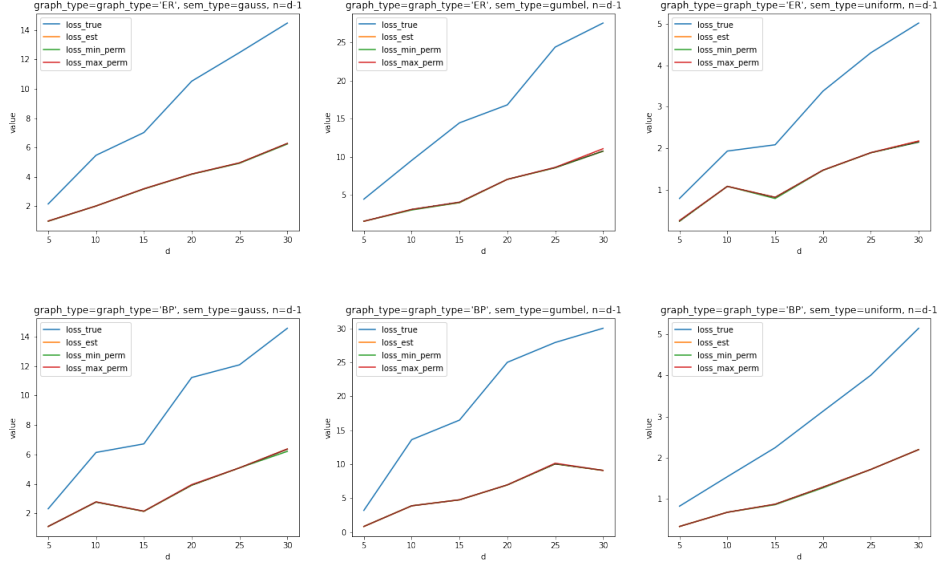


Figure 3: L-BFGS

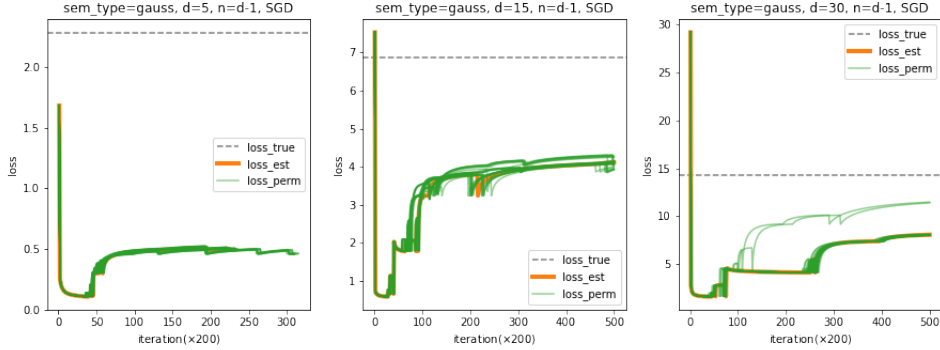


Figure 4: SGD

The nature question is that why NOTEARS is quite robust under permutations. Hence we use a three-node template to test if the algorithm could identify the actual edge direction under permutations and figure out how do the losses perform correspondingly. From Figure 5 and Figure 10, we could find that the algorithm does fail to identify the true direction but the L2 losses are optimized to the similarly same level under permutations.

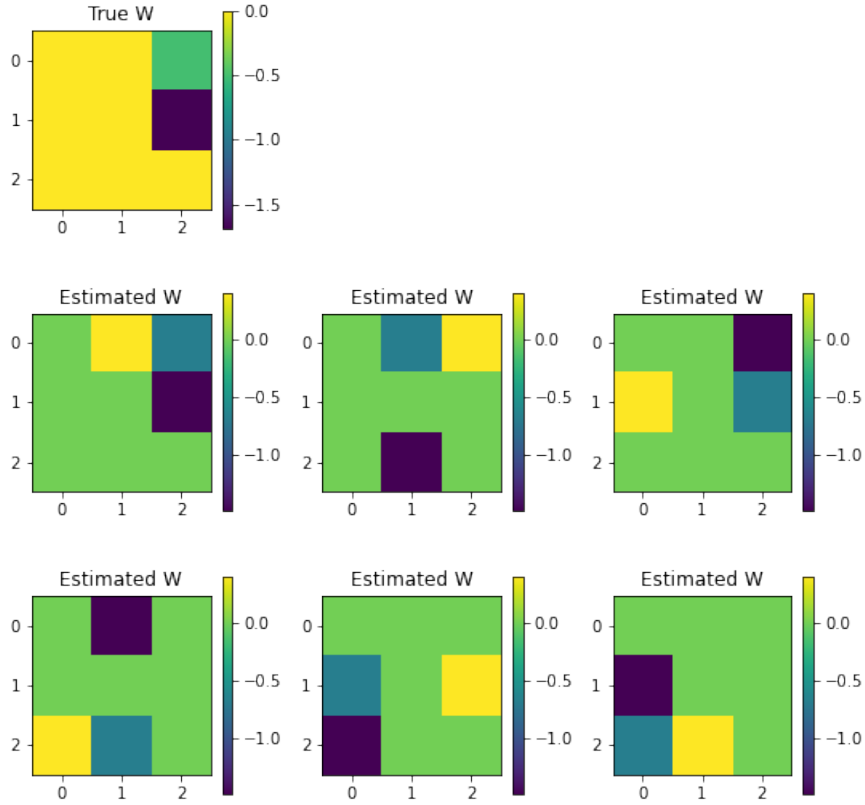


Figure 5: 3-node template

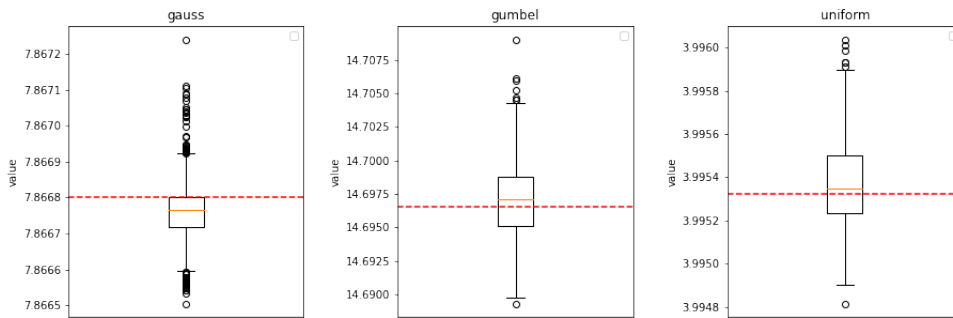


Figure 6: 3-node template repeated for 4 times with noise with variance of Gaussian distributed size (uniform one behaves similarly)

Actually, this phenomenon does not come forth only when we use some extreme examples, the algorithm fails to identify the true orientation even for constant variance and normal ER graphs as shown in Figure 7.

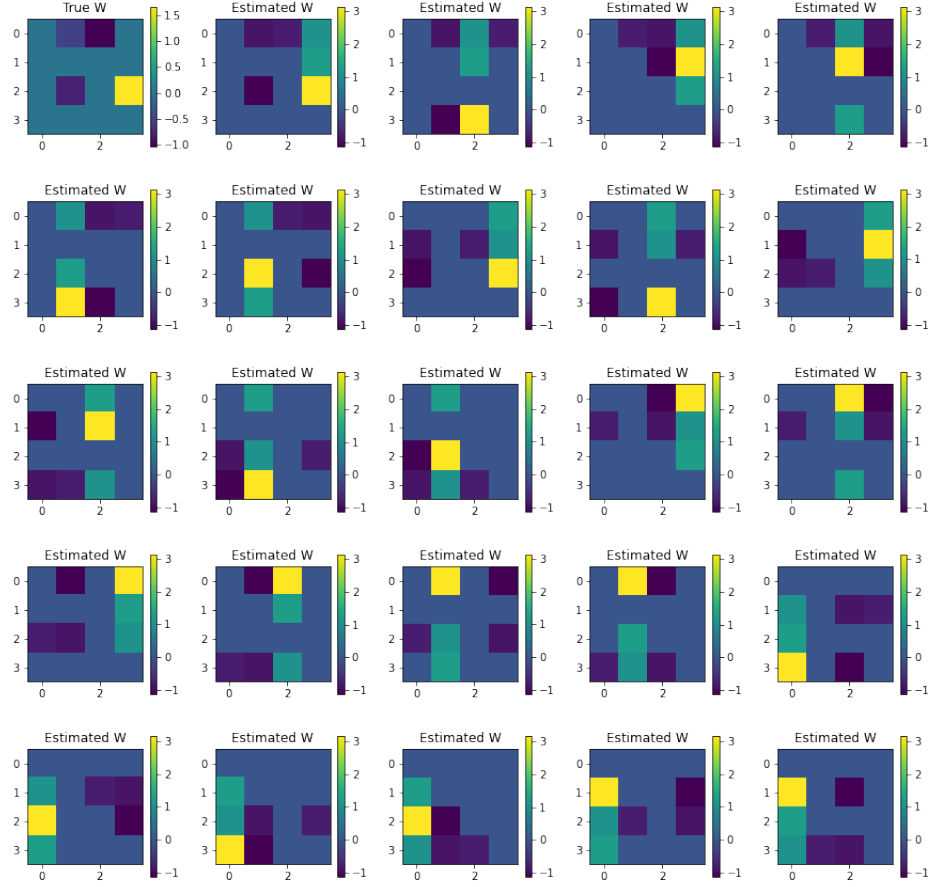


Figure 7: 4-node ER graph

Then I also explored if NOTEARS is robust towards data generated with non-constant variance and found that NOTEARS performs not so robust on data generated by noise with variance of Gaussian distributed size. The results are shown in Figure 8. However, previous papers have talked about this problem and Ignavier Ng, etc.<sup>[2]</sup> suggested using log-likelihood to avoid rescaling problem and proposed GOLEM. **However, we could find that GOLEM performs even worse than NOTEARS in terms of robustness under settings of heterogeneous variance from Figure 9.**



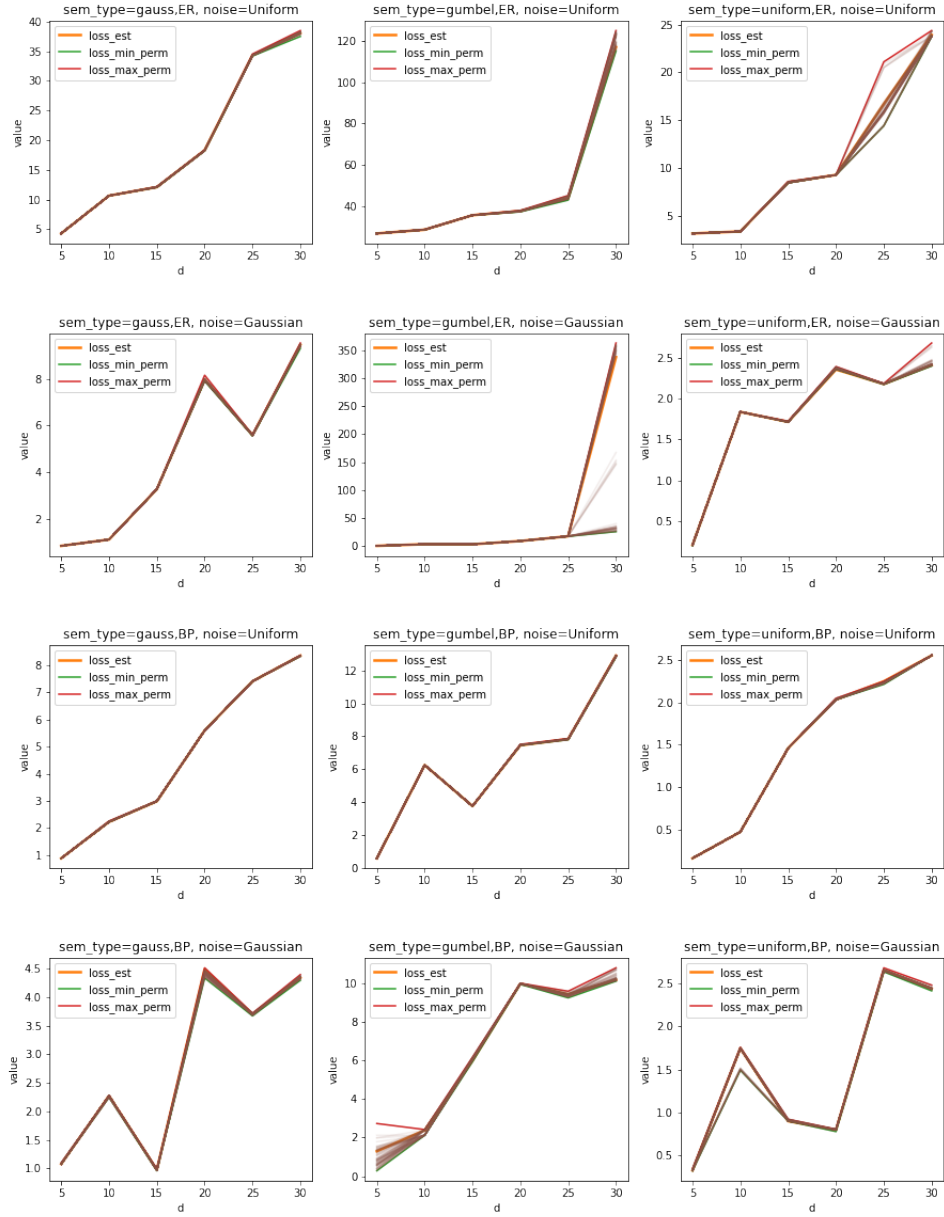


Figure 8: NOTEARS with non-constant variance

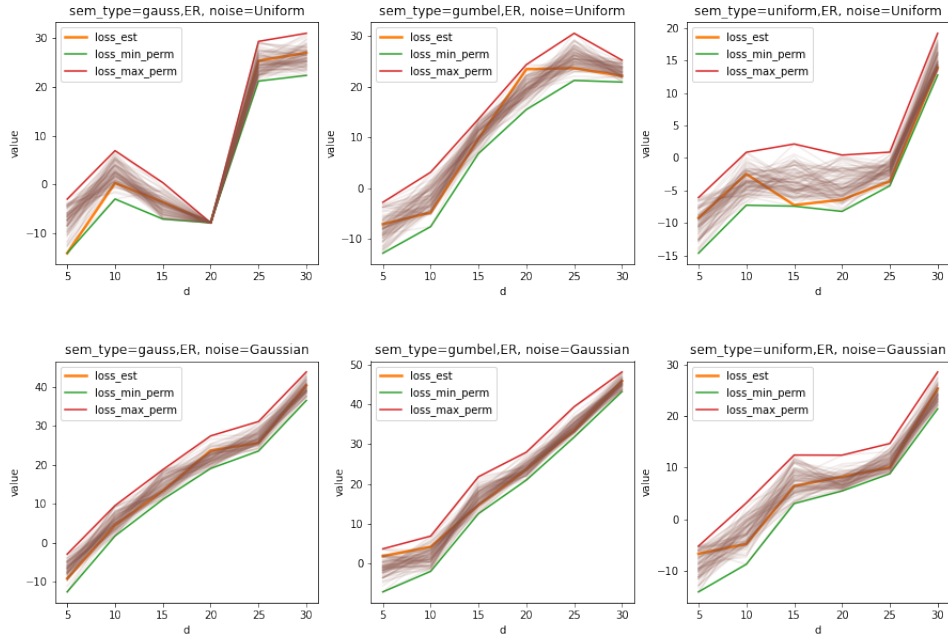


Figure 9: GOLEM with non-constant variance

As a supplement, DAGMA performs the same way in terms of failure in orientation identification and similar L2 loss after optimization under permutations as shown in Figure 10 and Figure 11. Meanwhile, it also shows signs of instability when variance is not constant in Figure 12.

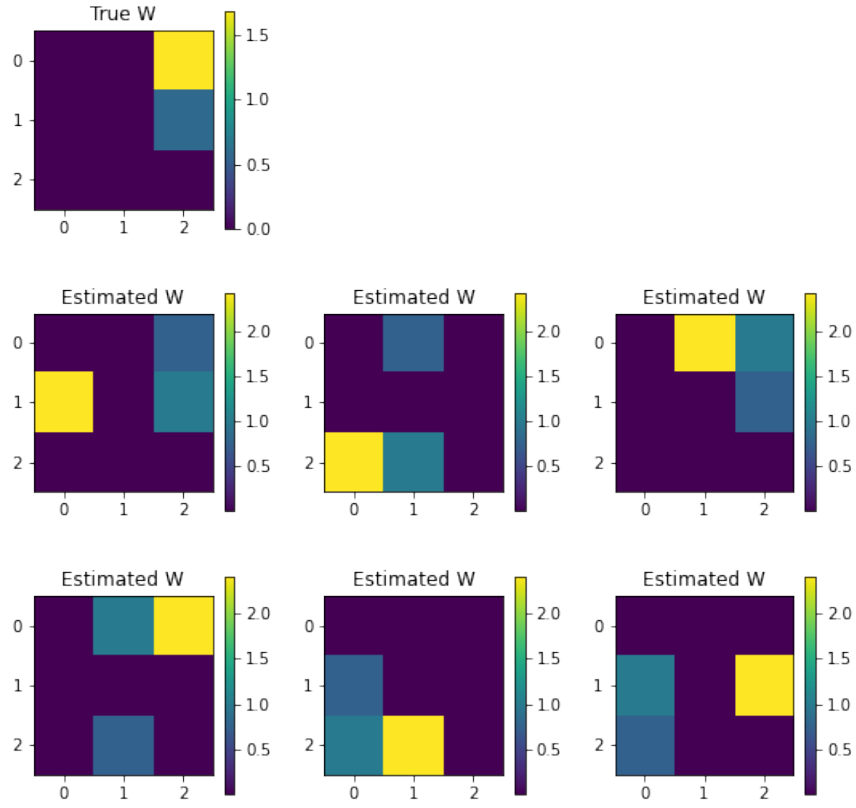


Figure 10: 3-node template by DAGMA

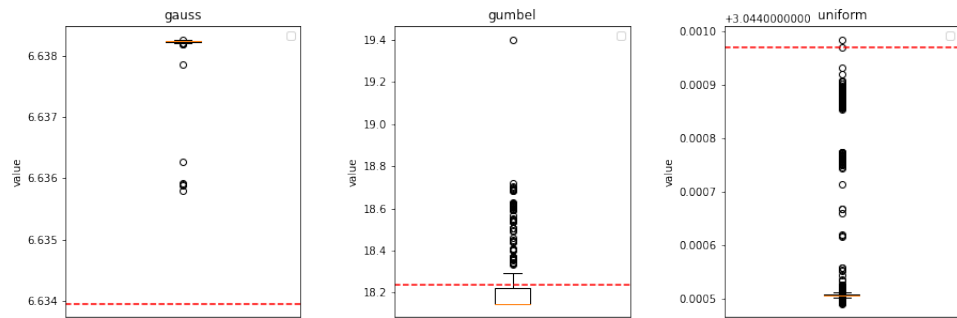


Figure 11: 3-node template repeated for 4 times with noise with variance of Gaussian distributed size (uniform one behaves similarly) by DAGMA

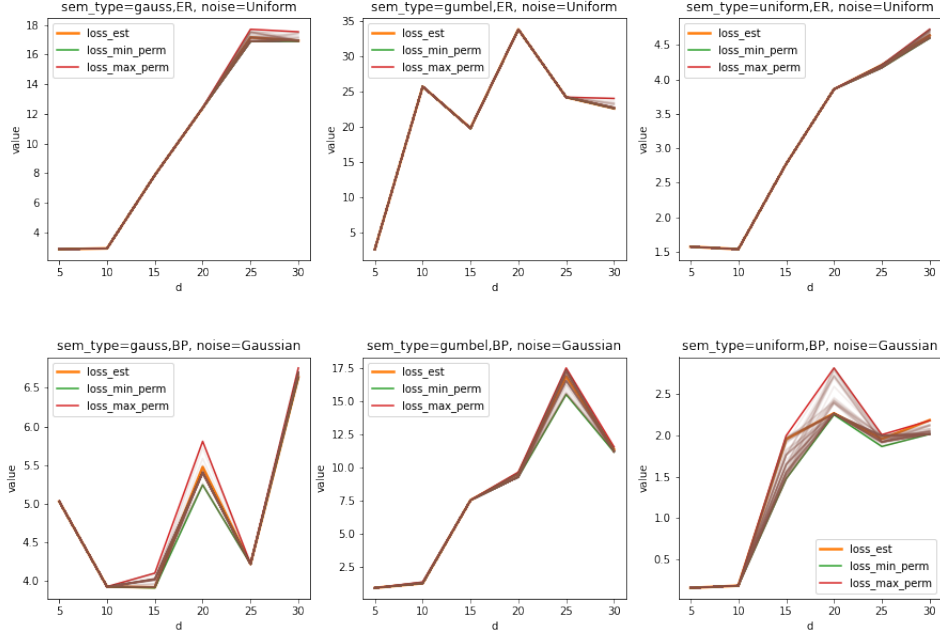


Figure 12: DAGMA with non-constant variance

Moreover, we wonder how does NOTEARS perform if data is discrete distributed. Surprisingly, from Figure 13, NOTEARS still perform well in discrete cases (data generated by logistic distribution) even when loss function is not a perfect match. It is worth noting that I did not introduce the heteroskedasticity problem here, because the general data generation function does not have relevant instructions, but perhaps this is an issue worth exploring.

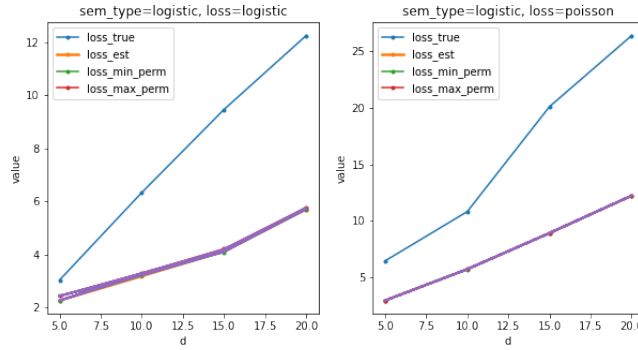


Figure 13: NOTEARS with data from logistic distribution

## 5 Limitations

First, notice the constraint function  $h(W) = \text{tr}(e^{W \circ W}) - d = \sum_{k=0}^{\infty} \frac{1}{k!} (W \circ W)^k - d$  is non-convex, leading to a non-convex objective function  $O(W, \rho)$ . This contributes to multiple minima and possibility of getting stuck in local minimum during optimization.

Second, the algorithm actually tells estimates from parsimonious linear SEMs rather than causal discoveries<sup>[3]</sup>. Sometimes the magnitude of estimates of coefficients in SEMs doesn't

imply the yes or no problem. Especially, the estimated  $W$  will be truncated by an arbitrary threshold.

Third, the algorithm is not robust to rescaling data. Marcus Kaiser and Maksim Sipos<sup>[3]</sup> show that first update step largely determines the final structure of estimated DAG. Also nodes with high variance are preferred to be sinks as opposed to sources. Moreover, simulated data are generated from additive noise models where child nodes have a higher variance than their parents<sup>[4]</sup>.

Later, Ignavier Ng, etc.<sup>[2]</sup> suggest using maximized log-likelihood (consider covariance matrix in the denominator to adapt to various data scales) and using soft DAG constraint (directly minimizing  $\ell + \lambda_1 \|W\|_1 + \lambda_2 h(W)$  without hard constraints).

## References

Xun Zheng, Bryon Aragam, Pradeep Ravikumar, and Eric P. Xing. Dags with no tears: Continuous optimization for structure learning. 2018.

Ignavier Ng, AmirEmad Ghassami, and Kun Zhang. On the role of sparsity and dag constraints for learning linear dags. 2021.

Marcus Kaiser and Maksim Sipos. Unsuitability of notears for causal graph discovery. 2021.

Alexander G. Reisach, Christof Seiler, and Sebastian Weichwald. Beware of the simulated dag! causal discovery benchmarks may be easy to game. 2021.