# DAGs with Tears: A Novel Structure Learning Method under Deep Learning Framework

Zhichao Chen, and Zhiqiang Ge, *Senior Member, IEEE*

*Abstract*—Bayesian network is a frequently-used method for fault detection and diagnosis in industrial processes. The basis of Bayesian network is structure learning which learns a directed acyclic graph from data. Since the search space in structure learning will scale super-exponentially with the increase of process variables, data-driven structure learning is a challenging problem. As a novel method for structure learning, DAGs with No Tears methods are being well studied in recent years due to their compatibility with deep learning framework. However, the DAGs with No Tears methods is far from application in industrial scenario due to problems in the gradient descent based solving stage and the post-processing stage. In this work, those problems are theoretically analyzed in detail by mathematical derivations. To solve these problems, the DAGs with Tears method is proposed by using mix-integer linear programming under the deep learning framework. In addition, prior knowledge is able to incorporate into the new proposed method, making structure learning more practical and useful in industrial processes. Finally, a numerical example and an industrial simulation example are adopted as case studies to demonstrate the superiority of the developed method.

*Index Terms*—Structure learning, directed acyclic graph, Bayesian network, gradient descent, mix-integer linear programming

## I. INTRODUCTION

**D**UE TO the high complexity mechanism and the lack of rigorous mathematical models, the mechanism-driven plant-level optimization is far from widespread application [1], [2]. Thanks to the great improvement of industrial intelligence and information technology, the acquisition and wide application of industrial big data have become possible. Therefore, the application of data-driven optimization is being a trending research topic. Quantities of new data-driven modeling method [3], [4] have recently been proposed especially using the deep learning method [5], which play a crucial part in the guidance of control, the increasing of economic value, and the guarantee of process safety. However, using the data acquired from the process directly could not promise the reliability of the model and results in the obstacle of the data-driven method promotion [6]. Meanwhile, as a branch of probabilistic graphical model, Bayesian network [7] is being more attractive, which can be regarded as a data structure that provides the skeleton for representing a joint distribution compactly in a factorized way [8]. As a compromise proposal of mechanism and data, Bayesian

network makes it possible to open the black box of data-driven models. The fault detection & diagnosis technology [9], [10] and soft sensor method [11] based on Bayesian network have been well studied and hence the research on the application of Bayesian network in industrial big data is of great importance.

The construction of Bayesian network consists of parameter learning and structure learning [8]. The structure learning is the basic of Bayesian network construction and a popular research topic in this field [12]. Actually, learning the directed acyclic graph (DAG) from the data directly as the graph of Bayesian network is an NP-Hard problem. The major difficulty is the combinational explosion of binary variables and the non-convexity for the acyclic constraint in the optimization problem. To learn the DAGs from data, three major methods are being used namely constraint based method, score based method, and DAGs with No Tears (NOTEAR) [13] method. Constraint based methods like PC algorithm [14] view a Bayesian network as a representation of independencies. They try to test for conditional dependence and independence in the data and then to find a network (or more precisely an equivalence class of networks) that best explains these dependencies and independencies. Score-based methods like K2 algorithm [15], MCMC algorithm [16], and Hill Climbing Search algorithm [17] view a Bayesian network as specifying a statistical model and then address learning as a model selection problem. Score based methods all operate on the same principle: Define a hypothesis space of potential models — and a scoring function that measures how well the model fits the observed data. The mentioned above methods are usually based on heuristic rule and could not consider all nodes simultaneously to find an optimal structure. Different from traditional heuristic methods, GoBNILP algorithm [18], as a score based method, used the mix-integer linear programming (MILP) to traverse all the nodes to maximize the score function. The prior knowledge can be added in the Bayesian network via the regulation of binary variables, while the combinational optimization and the corresponding constraints makes it to be an NP-Hard problem when solving the MILP linear programming problem. As a combination of score based and combinational optimization method, NOTEAR utilizes the statistical properties of the least square loss of structural equation model (SEM) in scoring DAGs. That is, the minimizer of the least square loss in the SEM provably recovers a true DAG with high probability on finite-samples and in high dimensions [19]. Therefore, the combinational optimization of DAG structure learning can be converted into continuous optimization problem and the computational efficiency and performance can be improved substantially.

The earliest NOTEAR method is a linear model, which makes it difficult to confront with the strong nonlinearity of process data. Therein, the least square loss of structural equation model in NOTEAR under deep learning framework is being studied in recent years to improve the accuracy when applying NOTEAR method to nonlinear data. Yu et al. [20] adopted the graph convolution operation to combine the VAE and NOTEAR as well as changed the constraint in NOTEAR simultaneously. The so-called DAG-GNN model is the first method that extends the NOTEAR to the nonlinear data under deep learning framework and the non-convexity constraint of the original is improved in a power form. Ng et al. [21] analyzed the SEM in an abstract function form and generalized the NOTEAR method to nonlinear case. Wang et al. [22] proposed a generative adversarial framework for NOTEAR to improve the DAG mining capability in nonlinear data.

Even though the NOTEAR has been well studied and achieved considerable improvement. In practice, the non-convexity of acyclic constraint is difficult to be satisfied under the deep learning framework whose parameters are updated via gradient descent method [23], which results in the infeasible solution in the final results constantly. Therein, the DAGs obtained from the NOTEAR methods should make post-processing to satisfy the acyclic constraint. Generally, the tear problem is converted to be a truncation problem [24], where some elements in the adjacent matrix will be set to 0 to avoid the tear problem. This post-processing method may have great perturbation of least square result. Meanwhile, to the best of the our knowledge, the NOTEAR is purely data-driven method, the prior knowledge cannot be added into the final DAG. This drawback reduces the reliability of the DAG since the models driven by mechanism and data may be more effective than those purely driven by data. To solve the two problems mentioned above, the DAGs with Tears (WITHTEAR) method is proposed in this work to take the advantages of NOTEAR and combinational optimization into account. The innovation of this work can be summarized as follows:

1) The occurrence of infeasible solution in NOTEAR method is theoretically analyzed under the gradient descent method, and the principle of the truncation operation in the post processing stage is analyzed in the perspective of least square perturbation.

2) The WITHTEAR method is formulated based on the NOTEAR method using MILP model. Besides, the MILP model considering prior knowledge of the data is combined into the DAG structure learning.

The paper is organized as follow: the problem statement is given in Section 2. In Section 3, the corresponding analyses of NOTEAR are provided in detail, and the WITHTEAR is given in Section 4. Two case studies are given in Section 5 to show the superiority of WITHTEAR baes on DAG-GNN model. Conclusions are drawn in Section 6.

## II. PROBLEM STATEMENT

The problem to be solved in this paper is given as follows: Givens are $n$ variables $x_i$ which belongs to a set $X$ defined as

$X = \{x_1, x_2, \ldots, x_n\}$. The prior knowledge of the variables is presented by the connection between variables using matrix $P$. The problem is how to combine NOTEAR and MILP model to learn Bayesian Network (DAG) under deep learning framework and perform the corresponding DAG topological structure using adjacent matrix $A$.

## III. THEORETICALLY ANALYSES OF NOTEAR METHOD

### A. The occurrence of infeasible solution

To better perform the WITHTEAR method, the analysis of NOTEAR method is given as follows. According to reference [13], the NOTEAR methods convert the traditional combinational optimization problem:

$$\min_A F(W)$$
$$s.t. G(W) \in DAGs$$

into a continuous programming problem shown as Problem P1:

$$(P1) \quad \min_{A,\Theta} \frac{1}{2n} \sum_{j=1}^{n} \left\| X_{(j)} - f(X_{(j)}|A,\Theta) \right\|_F^2 + \lambda \|A\|_1$$
$$s.t. h(A) = 0$$

The equality constraint can be written as Eq. (1) [13] or Eq. (2) [20]:

$$h(A) = Tr(e^{A \odot A}) - d = 0 \tag{1}$$

$$h(A) = Tr[(I + \gamma A \odot A)^d] - d = 0 \tag{2}$$

where the $Tr$ is the trace of matrix, $\odot$ is Hadmard product, $\gamma$ is hyper-parameter, $d$ is the dimension of square matrix $A$.

To solve this constrained optimization problem, according to the reference [25], the augmented Lagrangian method is adopted and the objective function of Problem P1 without constraint can be rewritten as shown in Eq. (3). The $\alpha$ and $\beta$ in Eq. (3) are Lagrange multiplier and penalty parameters respectively. When $c \to \infty$, the minimizer of Eq. (3) must satisfy $h(A) = 0$, in which case Eq. (3) is equal to the objective function of Problem P1.

$$\min_{A,\Theta} \frac{1}{2n} \sum_{j=1}^{n} \left\| X_{(j)} - f(X_{(j)}|A,\Theta) \right\|_F^2$$
$$+ \lambda \|A\|_1 + \alpha h(A) + \frac{\beta}{2} |h(A)|^2 \tag{3}$$

Hence the strategy is progressively increase $c$, and the Lagrange multiplier $\lambda$ is correspondingly updated as shown in Eq. (4) and (5).

$$\alpha_{k+1} = \alpha_k + \beta_k h(A_k) \tag{4}$$

$$\beta_{k+1} = \begin{cases} 10\beta_k, & if \, |h(A_k)| > 0.25 \, |h(A_{k-1})| \\ \beta_k, & otherwise \end{cases} \tag{5}$$

However, in the iteration process of the gradient descent method, the $A$ may not satisfy the acyclic constraint and results in an infeasible solution. To illustrate this phenomenon, the Problem P1 will be simplified as a linear form as Problem P2 shown. The constraint written in Eq. (1) and (2) will be discussed respectively.

$$(P2) \quad \min_A \frac{1}{2n} \sum_{j=1}^{n} \left\| X_{(j)} - X_{(j)}A \right\|_F^2$$
$$s.t. h(A) = 0$$

When Eq. (1) is adopted as constraint, the objective function can be formulated as Eq. (6) shown:

$$Loss = \frac{1}{2n} \sum_{j=1}^{n} \left\| X_{(j)} - X_{(j)}A \right\|_F^2 \tag{6}$$
$$+ \alpha(Tr(e^{A \odot A}) - d) + \frac{\beta}{2} \left| Tr(e^{A \odot A}) - d \right|^2$$

Assume that in the iteration of gradient descent, the $A_k$ calculated in the $k$-th time iteration satisfies the constraint, then the $A_{k+1}$ in the next time iteration is shown in Eq. (7).

$$A_{k+1} = A_k - LR \times \frac{\partial Loss}{\partial A_k} \tag{7}$$

where *LR* corresponds to the learning rate in the gradient descent method.

The Eq. (7) can be expanded [26], [27] in Eq. (8):

$$A_{k+1} = A_k - LR \times [\frac{1}{n} \sum_{j=1}^{n} X_{(j)}^T(X_{(j)}A_k - X_{(j)}) \tag{8}$$
$$+ 2\alpha A_k \odot (e^{A_k \odot A_k})^T]$$

then, Eq. (9) and (10) can be derived by multiplying $A_k$ and $A_{k+1}$ on both sides respectively.

$$A_{k+1} \odot A_k = A_k \odot A_k$$
$$- LR \times [\frac{1}{n} \sum_{j=1}^{n} X_{(j)}^T(X_{(j)}A_k - X_{(j)}) \tag{9}$$
$$+ 2\alpha A_k \odot (e^{A_k \odot A_k})^T] \odot A_k$$

$$A_{k+1} \odot A_{k+1} = A_k \odot A_{k+1}$$
$$- LR \times [\frac{1}{n} \sum_{j=1}^{n} X_{(j)}^T(X_{(j)}A_k - X_{(j)}) \tag{10}$$
$$+ 2\alpha A_k \odot (e^{A_k \odot A_k})^T] \odot A_{k+1}$$

Substitute Eq. (10) to Eq. (9), Eq. (11) can be given as follow:

$$A_k \odot A_k = A_{k+1} \odot A_{k+1}$$
$$+ LR \times [\frac{1}{n} \sum_{j=1}^{n} X_{(j)}^T(X_{(j)}A_k - X_{(j)}) \tag{11}$$
$$+ 2\alpha A_k \odot (e^{A_k \odot A_k})^T] \odot (A_{k+1} + A_k)$$

Finally, the LHS of Eq. (11) is part of constraint Eq. (1), Hence, Eq. (12) can be derived when substituting Eq. (11) to Eq. (1):

$$h\{A_{k+1} \odot A_{k+1}$$
$$+ LR \times [\frac{1}{n} \sum_{j=1}^{n} X_{(j)}^T(X_{(j)}A_k - X_{(j)}) \tag{12}$$
$$+ 2\alpha A_k \odot (e^{A_k \odot A_k})^T] \odot (A_{k+1} + A_k)\} = 0$$

Since the data is fed batch-by-batch, and the direction of gradient descent may not be determined due to the data shuffling every time of iteration. Besides, the *LR* as hyperparameters could not be appropriate to ensure the *A* be

feasible at the last time of iteration. Therein, Eq. (13) can't be permanent establishment during iteration of Eq. (12), which means that the algorithms will report an infeasible solution of *A* at the end of iteration.

$$[\frac{1}{n} \sum_{j=1}^{n} X_{(j)}^T(X_{(j)}A_k - X_{(j)}) + 2\alpha A_k \odot (e^{A_k \odot A_k})^T] \tag{13}$$
$$\odot (A_{k+1} + A_k) = 0$$

Next, Eq. (2) is adopted in Problem P2, and the objective function can be written as Eq. (14):

$$Loss = \frac{1}{2n} \sum_{j=1}^{n} \left\| X_{(j)} - X_{(j)}A \right\|_F^2$$
$$+ \alpha\{Tr[(1 + \gamma A \odot A)^m] - d\} \tag{14}$$
$$+ \frac{\beta}{2} \left| Tr[(1 + \gamma A \odot A)^m] - d \right|^2$$

Similarly, the $A_{k+1}$ can be derived as shown in Eq. (15):

$$A_{k+1} = A_k - LR \times \{\frac{1}{n} \sum_{j=1}^{n} X_{(j)}^T(X_{(j)}A_k - X_{(j)})$$
$$+ 2\alpha A_k \odot \sum_{k=1}^{d} C_d^k k \left[\gamma(A_k \odot A_k)^{k-1}\right]^T\} \tag{15}$$

where the $C$ stands for the combinationtorial number and after suitable transformation, Eq. (15) can be rewritten as Eq. (16):

$$A_k \odot A_k = A_{k+1} \odot A_{k+1}$$
$$+ LR \times \{\frac{1}{n} \sum_{j=1}^{n} X_{(j)}^T(X_{(j)}A_k - X_{(j)})$$
$$+ 2\alpha A_k \odot \sum_{k=1}^{d} C_d^k k \left[\gamma(A_k \odot A_k)^{k-1}\right]^T\} \tag{16}$$
$$\odot (A_k + A_{k+1})$$

Substituting Eq. (16) to Eq. (2), the Eq. (17) similar to Eq. (12) can be derived which could not promise a feasible solution after iteration.

$$h[A_{k+1} \odot A_{k+1}$$
$$+ LR \times \{\frac{1}{n} \sum_{j=1}^{n} X_{(j)}^T(X_{(j)}A_k - X_{(j)})$$
$$+ 2\alpha A_k \odot \sum_{k=1}^{d} C_d^k k \left[\gamma(A_k \odot A_k)^{k-1}\right]^T\} \tag{17}$$
$$\odot (A_k + A_{k+1})] = 0$$

Last, the nonlinear form is adopted in Problem P1 and the similar equations can be given in Eq. (18) and (19), which the similar phenomena can be derived.

$$h\{A_{k+1} \odot A_{k+1} + LR$$
$$\times [\frac{1}{n} \sum_{j=1}^{n} (\frac{\partial f(X_{(j)}|A_k, \Theta)}{\partial A_k})^T (f(X_{(j)}|A_k, \Theta) - X_{(j)})$$
$$+ 2\alpha A_k \odot (e^{A_k \odot A_k})^T] \tag{18}$$
$$\odot (A_{k+1} + A_k)\} = 0$$

$$h\{A_{k+1} \odot A_{k+1} + LR$$
$$\times [\frac{1}{n}\sum_{j=1}^{n}(\frac{\partial f(X_{(j)}|A_k,\Theta)}{\partial A_k})^T (f(X_{(j)}|A_k,\Theta) - X_{(j)})$$
$$+ 2\alpha A_k \odot \sum_{k=1}^{d} C_d^k k \left[\gamma(A_k \odot A_k)^{k-1}\right]^T]$$
$$\odot (A_{k+1} + A_k)\} = 0 \tag{19}$$

It should be noticed that, during the practice, if the L1 norm regularization is added in the objective function, the $A$ will tends to be 0. The reason can be interpreted through observing Eq. (12), (17), (18) and (19). With the increase of $\gamma$ as the iterations increases, the existence of the factor shown in the LHS of Eq. (20) gained from the above-mentioned equations tends to be 0, meanwhile, the L1 norm will forced the $A$ to be sparsity [28], and hence the $A$ will tend to be 0 in the iteration.

$$A_{k+1} + A_k = 0 \tag{20}$$

In conclusion, due to the limit of gradient descent method and the non-convexity of the equation constraint, the DAGs acquired via NOTEAR methods will not promise to be acyclic. Therein, the post-processing operation accompanied by the NOTEAR methods arises. In the next section, the principle of post-processing operation (also known as truncate as mentioned before) will be stated for better introducing the WITHTEAR method.

### B. The principle of the post-processing operation

In post-processing operation, a threshold will be set and the elements in matrix $A$ lower than this threshold will be set to 0, and the new matrix $A$ is obtained. If the new matrix $A$ satisfies the acyclic condition, then the new matrix $A$ will be output as the final results, otherwise the new matrix $A$ is cyclic, then the threshold will increase and the matrix will be truncated again until it is acyclic. In the truncate process, the elements of $A$ will make the optimal least square result deviate from the optimal point. Therefore, the rationality of the truncation operation should be analyzed, which to the best of our knowledge has not been studied before. A simple analysis of the least square will be provided as follows to show the principle of the post-processing operation.

Firstly, the disturbance of the linear least square problem will be considered. According to the least square problem shown as Eq. (21):

$$\tilde{X} = XA \tag{21}$$

If a perturbation $\delta A$ is added to $A$, the regression perturbation on the LHS of Eq. (21) can be given as follow:

$$\delta \tilde{X} = X \times \delta A \tag{22}$$

Then, the inequality can be derived:

$$\left\|\delta \tilde{X}\right\|_2 \leqslant \|\delta A\|_2 \|X\|_2 \tag{23}$$

Both sides of Eq. (23) can be divided by the L2 norm of matrix $X$ and the relative error of the perturbation can be given

as Eq. (24). Similarly, for the non-linear form, the perturbation is given as Eq. (25).

$$\frac{\left\|\delta \tilde{X}\right\|_2}{\left\|\tilde{X}\right\|_2} \leq \frac{\|A\|_2 \|X\|_2}{\|XA\|_2} \frac{\|\delta A\|_2}{\|A\|_2}$$
$$= (\frac{\|A\|_2 \|X\|_2}{\|XA\|_2 \times \|A\|_2}) \times \|\delta A\|_2 \tag{24}$$

$$\frac{\left\|\delta \tilde{X}\right\|_2}{\left\|\tilde{X}\right\|_2} \leq \frac{\left\|\frac{\partial f(X|A,\Theta)}{\partial A}\right\|_2}{\|f(X|A,\Theta)\|_2} \|\delta A\|_2 \tag{25}$$

Summarizing Eq. (24) and (25), it can be concluded that, once the matrix $A$ and paramters of regression model is determined, the relative error merely depends on the perturbation magnitude, which means that, to make the deterioration of the least square result as less as possible, the elements of $A$ will be set to be 0 from small to large until the graph constructed by $A$ is acyclic. This process can adopt tear operation to fulfill. However, for convenience, the NOTEAR, makes the tear problem into a truncation problem as mentioned before. Note that, the truncation operation may deteriorate the least square loss greater than that of the tear operation. Meanwhile, the prior knowledge can't merge with the knowledge learned from data in NOTEAR through roughly truncation. Therefore, the following WITHTEAR method will be proposed to solve the two problems mentioned above.

## IV. DAGs with Tears Method

### A. Loops Tear by MILP problem

In the previous section, to promise the minimum perturbation of the least square problem and tear all loops of the graph simultaneously, the number of elements in $A$ should be changed as few as possible. Therein, MILP model can be formulated to fulfill the two goals as stated above and the corresponding method is named DAGs with Tears (WITHTEAR). Before showing the WITHTEAR method, the following concept will be stated as follow for a better understanding of WITHTEAR method.

Loop matrix given in Eq. (26) is a matrix with connection between node (also known as stream, whose set is written as STR) as column and loop $i$ as row, $[u_{i,j}]$. If a loop $i$ includes a stream $j$, the element in the loop matrix $u_{i,j} = 1$, otherwise $u_{i,j} = 0$.

$$U = \begin{bmatrix} 0 & \cdots & 1 \\ \vdots & \ddots & \vdots \\ 1 & \cdots & 0 \end{bmatrix} \tag{26}$$

Based on Eq. (26), innovated by [29], [30], the loop tearing cost is introduced as Eq. (27) to measure the cost of breaking the stream.

$$Cost = \sum_{j \in STR} w_j \times y_j, y_j \in \{0, 1\} \tag{27}$$

where $w_j$ is the weight of the stream, which can be converted from the coefficient of $A$, and $y_j$ is binary variable to be solved. If $y_j$ is 0, then the stream should be tear, vice versa. To tear

all loops, the corresponding constraints is given in Eq. (28). This constraint indicates that the loop in matrix $U$ should be tore at least. By adopting this constraint, the acyclic of the graph formulated from matrix $A$ can be promised.

$$\sum_{j \in STR} u_{ij} y_j \geq 1, u \in U, y_j \in \{0, 1\} \qquad (28)$$

In all, the MILP problem can be given in problem P3. By solving problem P3, the loops of matrix $A$ can be tore and the final DAG can be form with the minimum deterioration of least square result.

$$(P3) \quad \min Cost = \sum_{j \in STR} w_j \times y_j$$
$$s.t. \begin{cases} \sum_{j \in STR} u_{ij} y_j \geq 1, u \in U \\ y_j \in \{0, 1\} \end{cases}$$

### B. Loops Tear combing prior knowledge and MILP problem

Note that, the prior knowledge is not added in the MILP problem P3. Hence, in this section, Problem P3 will be extended to combine prior knowledge with the addition of logical propositions and the corresponding disjunctions. The prior knowledge can be given in 3 scenarios:(1), the existence of stream $j$ is unknown; (2), the existence of stream $j$ is obligatory; (3), the existence of stream $j$ is forbidden. For scenario (1), the corresponding elements in matrix A can be set to be 0 before the solving of MILP problem P3. Therein, the scenario (2) and (3) will be discussed as follow: To show the relative position relationships, the logical variables [31] are defined as below:

1). $V_{1,j}$ If the existence of stream $j$ is unknown, $V_{1,j}$ is True.

2). $V_{2,j}$ If the existence of stream $j$ is obligatory, $V_{2,j}$ is False.

The disjunctions can be given as Eq. (29) shown:

$$\begin{bmatrix} V_{1,j} \\ UB_{1,j} = 1.0 \\ LB_{1,j} = 0.0 \\ j \in STR \end{bmatrix} \vee \begin{bmatrix} V_{2,j} \\ UB_{2,j} = 0.5 \\ LB_{2,j} = 0.0 \\ j \in STR \end{bmatrix} \qquad (29)$$
$$V_{1,j}, V_{2,j} \in \{True, False\}$$

where $UB$ and $LB$ are the upper bound and lower bound of binary variable $y_j$ defined in Eq. (28). The extra constraint for the binary variable $y_j$ is given as Eq. (30):

$$LB_j \leq y_j \leq UB_j \qquad (30)$$

there are two constraints in Eq. (29), while only one constraint will take effect, and the remains do not work. Therefore, each logical variable should satisfy the constraint shown in Eq. (30):

$$V_{1,j} \veebar V_{2,j} \qquad (31)$$

Finally, the optimization problem can be formulated as problem P4, and the value of the logical variable can be derived from matrix $P$ that contains prior knowledge.

$$\min Cost = \sum_{j \in STR} w_j \times y_j$$
$$(P4) \quad s.t. \begin{cases} \sum_{j \in STR} u_{ij} y_j \geqslant 1, \ u \in U \\ y_j \in \{0, 1\} \\ LB_j \leqslant y_j \leqslant UB_j, \ j \in STR \\ \begin{bmatrix} V_{1,j} \\ UB_{1,j} = 1.0 \\ LB_{1,j} = 0.0 \\ j \in STR \end{bmatrix} \vee \begin{bmatrix} V_{2,j} \\ UB_{2,j} = 0.5 \\ LB_{2,j} = 0.0 \\ j \in STR \end{bmatrix} \\ V_{1,j}, V_{2,j} \in \{True, False\} \\ V_{1,j} \veebar V_{2,j} \end{cases}$$

### C. The algorithm for DAGs with Tears method

---
**Algorithm 1:** Pseudo-code of WITHTEAR Method

---

**input** : Data $X$, Parameter $\beta_{max}$, the maximum iteration time of deep learning model *Epoch*, the initial value of $\alpha$, $\beta$, and matrix $A$. The coefficient of L1 norm $\lambda$. The initial value of the best objective function of Problem P1 $Loss_{best} = \infty$, and matrix $A$ $A_{best} =$ None. The matrix which loads the prior knowledge $P$. The hyper-parameter $\omega$

**output**: The best matrix $A_{best}$

1 **Training Stage**
2 **while** $\beta \leqslant \beta_{max}$ **do**
3    **for** $i = 1 : Epoch$ **do**
4       
$$A_{k+1} = \begin{array}{c} \arg\min_A \frac{1}{2n} \sum_{j=1}^{n} \left\| X_{(j)} - f(X_{(j)}|A_k, \Theta) \right\|_F^2 \\ + \lambda \|A\|_1 + \alpha h(A_k) + \frac{\beta}{2} |h(A_k)|^2 \end{array}$$
      **if** $\frac{\sum_{j=1}^{n} \left\| X_{(j)} - f(X_{(j)}|A_k, \Theta) \right\|_F^2}{2n} \leqslant Loss_{best}$ **then**
5          $Loss_{best} = \frac{\sum_{j=1}^{n} \left\| X_{(j)} - f(X_{(j)}|A_k, \Theta) \right\|_F^2}{2n}$ ;
6          $A_{best} = A_k$ ;
7    $\alpha_{k+1} = \alpha_k + \beta_k h(A_k)$ ;
8    $\beta_{k+1} = \begin{cases} 10 \times \beta_k, \ if \ |h(A_k)| > 0.25 |h(A_{k-1})| \\ \beta_k, \ otherwise \end{cases}$ ;

---

According to Problem P3 and P4, the pseudo-code of WITHTEAR method can be given based on the results of NOTEAR as Algorithm 1. From the pseudo-code, it can be concluded that the training stage has no differences between NOTEAR method, any NOTEAR under deep learning framework can be used to acquire matrix $A$.

In the tear stage, the post processing strategy is changed comparing to the truncation operation adopted in NOTEAR

method. By solving the MILP problem, the prior knowledge of data can be added into the DAGs structure learning increasing the validity of the DAG structure mined from data. Meanwhile, comparing to traditional DAG structure learning using MILP like GoBNILP, the binary variables in this research merely depend on the number of streams rather than that of the parent node, results in a reduction in computation complexity.

---

**9 Tear Stage**

10     **Preprocessing**    For the elements in A lower than the hyper-parameter $\omega$ and the existence of corresponding streams are forbidden in $P$, the 0 should be placed in the corresponding location. For those should streams be existence in $P$ while the corresponding elements are 0 in $A$, the corresponding elements set to be $\|A\|_{\max}$ ;

11 **while** *Loops exist in matrix* $A_{best}$ **do**

12      Formulate loop matrix $U$;

13      Solve Problem $(P3)$ or $(P4)$;

14      For those streams should be tore, set the corresponding elements in $A_{best}$ to 0;

15      Detect the existence of loop in A

---

## V. CASE STUDIES

### A. Numerical example

In this section, the DAG-GNN without L1 norm will be adopted to test the effectiveness WITHTEAR method. More Detailed of DAG-GNN model [20] are listed in the appendix. The nonlinear data simulation shown in Eq. (31) from [13] is modified as the numerical example.

$$X = \tanh(XW^T) + \cos(XW^T) + \sin(XW^T) + z, z \sim \mathcal{N}(0, I) \tag{32}$$

The matrix $W$ is upper triangular matrix, and the prior knowledge is the existence of streams connected by the nodes in lower triangular place of matrix $A$ is forbidden. A total of 5, 000 data are used for the structure learning.

The score functions are false discovery rate (FDR), true positive rate (TPR), false positive rate (FPR), and structure hamming distance (SHD). The corresponding expression are given in Eq. (32) $\sim$ (35).

$$FDR = \frac{R + FP}{TEE} \tag{33}$$

$$TPR = \frac{TP}{T} \tag{34}$$

$$FPR = \frac{R + FP}{F} \tag{35}$$

$$SHD = E + M + R \tag{36}$$

In Eq. (32), the *R* and *FP* are edges predicted reversed and not in the undirected skeleton of the true graph, respectively. The *TEE* is the total number of estimated edges. In Eq. (33), *TP* and *T* are correct direction edges and true edges respectively.

The *F* is the non-edges in the ground truth graph in Eq. (34). The *E* is the extra edges from the skeleton while the *M* is the missing edges from the skeleton.

Fig. 1 proposed the predicting results and table 2 indicates
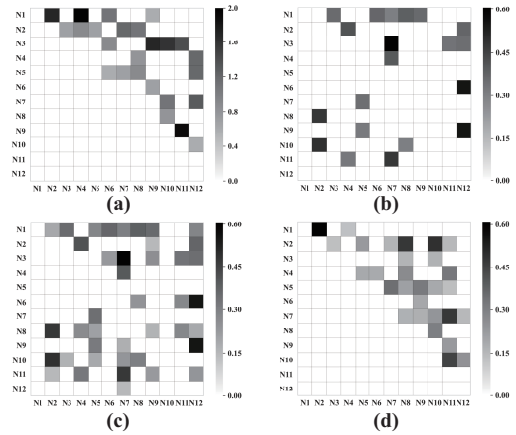


Fig. 1. The adjacent matrix (a), ground truth; (b), DAG-GNN; (c), DAGs with Tears (Problem P3); (d), DAGs with Tears (Problem P4).

the score function corresponds to the graphs in Fig. 1. By using WITHTEAR method, the *FDR*, *SHD* and *FPR* are 5.1∼24.0%, 14.6∼46.3%, 46.9∼59.4% lower than that of the DAG-GNN, respectively. If the prior is knowledge is not given in WITHTEAR, the *TPR* will 66.67 % lower than that of the DAG-GNN. While the situation is changed thanks to the prior knowledge matrix *P* is adopted in the WITHTEAR and results in the *TPR* 77.78% higher than that of the DAG-GNN. It can be seen that the prior knowledge on the coefficient matrix *W* can help the elimination of the false edge so as to increase the *TPR*.

TABLE I
THE SCORE FUNCTIONS OF THE CORRESPONDING METHODS

| ID | FDR | TPR | SHD | FPR |
|---|---|---|---|---|
| DAG-GNN | 0.81 | 0.36 | 41 | 0.78 |
| WITHTEARS(P3) | 0.77 | 0.12 | 35 | 0.41 |
| WITHTEARS(P4) | 0.62 | 0.64 | 22 | 0.32 |

### B. Tennessee-Eastman Process

As shown in Fig. 2, the Tennessee Eastman (TE) process has been widely used as a benchmark process [32] to compare evaluated process monitoring and fault diagnosis algorithms. This process includes five major units: a condenser, a reactor, a stripper, a recycle compressor, and a vapor/liquid separator. There are 41 measured variables, 12 manipulated variables. In this study, 33 variables listed in Table 3 are selected to demonstrate the effectiveness of the DAGs with Tears method. The dataset used for modeling contains a total of 1, 460 data.

Due to the ground-truth DAG structure is unknown in TE process, the BGe score [33] and Gaussian BIC score [34] are adopted to score the DAGs [35]. The matrix P which contains prior knowledge is given in Fig. 3.
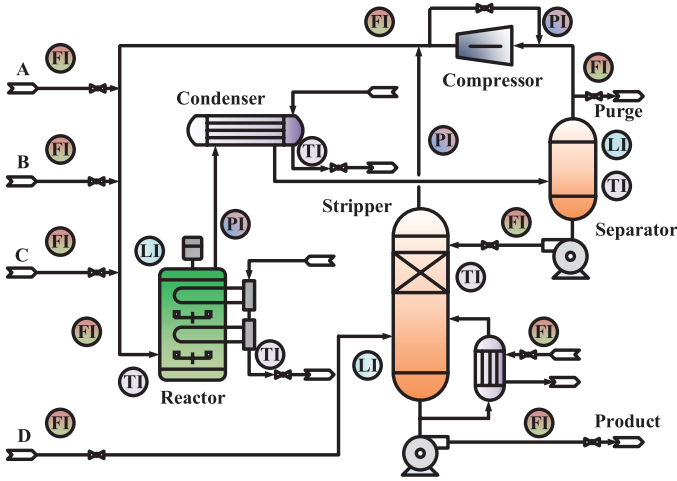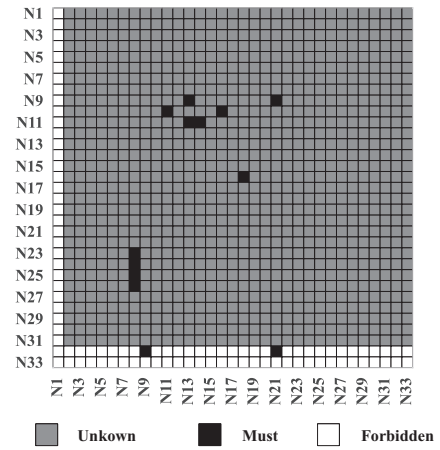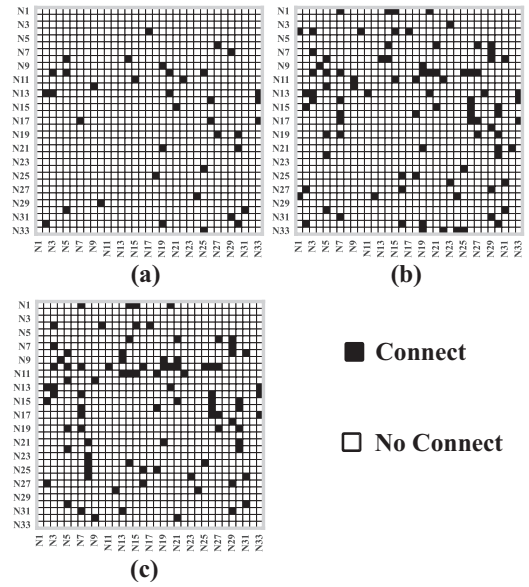
Fig. 2. The flowsheet of Tennessee-Eastman process.



Fig. 3. The prior knowledge matrix $P$.



Fig. 4. The adjacent matrix $A$ (a), DAG-GNN; (b), WITHTEAR (P3); (c), WITHTEAR (P4)

TABLE II
THE VARIABLES TO LEARN DAG STRUCTURE IN TE PROCESS

| No. | Measured Variables | No. | Measured Variables |
|---|---|---|---|
| N1 | A feed | N18 | Stripper temperature |
| N2 | D feed | N19 | Stripper steam flow |
| N3 | E feed | N20 | Compressor work |
| N4 | Total feed | N21 | Reactor cooling water outlet temperature |
| N5 | Recycle flow | N22 | Separator cooling water outlet temperature |
| N6 | Reactor feed rate | N23 | D feed flow |
| N7 | Reactor Pressure | N24 | E feed flow |
| N8 | Reactor level | N25 | A feed flow |
| N9 | Reactor temperature | N26 | Total feed flow |
| N10 | Purge rate | N27 | Compressor recycle valve |
| N11 | Product separator temperature | N28 | Purge valve |
| N12 | Separator level | N29 | Separator product liquid flow |
| N13 | Separator pressure | N30 | Stripper product liquid flow |
| N14 | Separator underflow | N31 | Stripper steam valve |
| N15 | Stripper level | N32 | Reactor cooling water flow |
| N16 | Stripper pressure | N33 | Condenser cooling water flow |
| N17 | Stripper underflow | | |

The DAG-GNN, WITHTEAR (Problem P3), and WITH-TEAR (Problem P4) are given in Fig. 4. The corresponding scores are listed in Table 4. The BGe score using WITHTEAR is 29.7 ∼ 34.8 % higher than that of the DAG-GNN, and the Gaussian BIC score is 27.8 ∼ 28.6 % higher than that of the using DAG-GNN. The behavior of WITHTEAR on the score function indicates that using the MILP problem to tear the matrix rather than roughly truncate will result in a better graph that represents the data. Note that, when solving problem P3 and P4, the BGe score and Gaussian BIC score are different, which indicates that the prior knowledge may

TABLE III
THE SCORES OF THE CORRESPONDING METHODS

| ID | BGe | Gaussian-BIC |
|---|---|---|
| DAG-GNN | -13, 1232.0 | -52, 982.1 |
| WITHTEARS (P3) | -85, 537.5 | -38, 237.4 |
| WITHTEARS (P4) | -92, 276.7 | -37, 838.3 |

exist bias and hence the prior should be given prudently when using WITHTEAR.

## VI. CONCLUSIONS

This paper proposed the WITHTEAR to improve the problem when applying NOTEAR methods to learn DAG structure under the deep learning framework. Firstly, the reason for the existence of infeasible solution results from the acyclic constraints adopted in the NOTEAR methods has been analyzed solved by the gradient descent method. After that, the principle for the truncation in post-processing operation was stated from the perspective of least square perturbation analysis. Based on this principle, the MILP models were formulated considering the tear cost and prior knowledge simultaneously in the final DAG formation. Finally, two case studies were carried out to demonstrate the effectiveness of the WITHTEAR method comparing to NOTEAR method using DAG-GNN model as baseline. The concentration of future work may be on the

hyper-parameter before the MILP problem solving to reduce the burden of loop detection algorithm or the directly tear method face to the adjacent matrix.

## APPENDIX A
### THE DERIVATION OF MATRIX DERIVATIVE

The derivation of matrix derivative in Section III will be stated as follow. First, consider Eq. (1) as constraint, define $\phi$ as Eq. (A.1)shown:

$$\phi = Tr(e^{A \odot A}) = Tr(I + \sum_{k=1}^{\infty} \frac{(A \odot A)^k}{k!}) \qquad (A.1)$$

Then, the matrix derivate of trace function [26], [27] can be given as Eq. (A.2):

$$d\phi = \sum_{k=0}^{\infty} \frac{\left[(A \odot A)^k\right]^T}{k!} : d(A \odot A)$$

$$= \sum_{k=0}^{\infty} \frac{\left[(A \odot A)^k\right]^T}{k!} : (dA \odot A + A \odot dA)$$

$$= \sum_{k=0}^{\infty} \frac{\left[(A \odot A)^k\right]^T}{k!} : 2A \odot dA \qquad (A.2)$$

where the : is the Frobenius / trace (inner) product. Then the Hardmard product and the Frobenius product can be commuted and the Eq. (A.3) can be derived

$$d\phi = \sum_{k=0}^{\infty} \frac{\left[(A \odot A)^k\right]^T}{k!} \odot 2A : dA \qquad (A.3)$$

Hence, the derivate of Eq. (1) to matrix $A_k$ is shown as Eq. (A.4)

$$\frac{\partial \phi}{\partial A} = 2A \odot \sum_{k=0}^{\infty} \frac{\left[(A \odot A)^k\right]^T}{k!} = 2A \odot (e^{A \odot A})^T \qquad (A.4)$$

Similarly, the matrix derivate of Eq. (2) is given as follow, define $\zeta$ as Eq. (A.5) shown.

$$\zeta = Tr[(I + \gamma A \odot A)^d]$$

$$= Tr[\sum_{k=1}^{d} C_d^k (I)^{d-k} (\gamma A \odot A)^k]$$

$$= Tr[\sum_{k=1}^{d} C_d^k (\gamma A \odot A)^k] \qquad (A.5)$$

Then, Eq. (A.6) can be derived:

$$d\zeta = \sum_{k=1}^{d} C_d^k k \left[(\gamma A \odot A)^{k-1}\right]^T : d(A \odot A)$$

$$= \sum_{k=1}^{d} C_d^k k \left[(\gamma A \odot A)^{k-1}\right]^T : 2A \odot dA \qquad (A.6)$$

Commute the Hardmard product and Frobenius product shown as Eq. (A.7).

$$d\zeta = 2A \odot \sum_{k=1}^{d} C_d^k k \left[(\gamma A \odot A)^{k-1}\right]^T : dA \qquad (A.7)$$

Finally, the matrix derivate of Eq. (2) can be given as follow:

$$\frac{\partial \zeta}{\partial A} = 2A \odot \sum_{k=1}^{d} C_d^k k \left[(\gamma A \odot A)^{k-1}\right]^T \qquad (A.8)$$

## APPENDIX B
### THE DAG-GNN MODEL

In Appendix B, the DAG-GNN model will be introduced. The architecture of DAG-GNN model is shown in Fig. B.1.
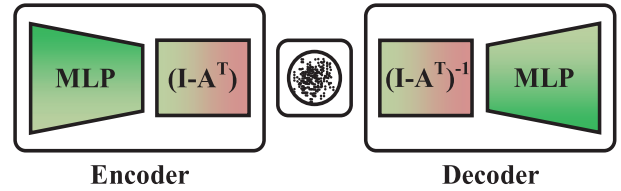


Fig. B.1.  The architecture of DAG-GNN model

The expression of encoder and decoder are given in Eq. (B.1) and (B.2) [20], respectively:

$$[M_Z | \log S_Z] = (I - A^T) MLP(X) \qquad (B.1)$$

$$[M_X | \log S_X] = MLP((I - A^T)^{-1} Z) \qquad (B.2)$$

where the $M$ and $S$ are mean and variance respectively. The *MLP* is the multi-layer perceptron and defined as Eq. (B.3). The $W_1$, $W_2$, $b_1$, and $b_2$ in Eq. (B.3) are parameters to be learnt via gradient descent.

$$MLP(X) = \text{ReLU}(X W_1^T + b_1) W_2^T + b_2 \qquad (B.3)$$

Therefore, the loss function can be given as Eq. (B.4) where the $D_{KL}$ is the Kullback-Leibler divergence, and the $p(Z)$ is the prior distribution.

$$Loss = E_{q(Z|X)}[\log p(X|Z)] + D_{KL}[q(Z|X)||p(Z)] \quad (B.4)$$

The RHS of Eq. (B.4) can be given as follow, where $m$ is theencoder dimension, $d$ is the number of variables, $L$ is the number of variable, and $c$ is a constant:

$$D_{KL}[q(Z|X)||p(Z)]$$

$$= \frac{1}{2} \sum_{i=1}^{m} \sum_{j=1}^{d} (S_Z)_{i,j}^2 + (M_Z)_{i,j}^2 - 2 \log (S_Z)_{i,j} - 1 \quad (B.5)$$

$$E_{q(Z|X)}[\log p(X|Z)]$$

$$\approx -\frac{1}{L} \sum_{l=1}^{L} \sum_{i=1}^{m} \sum_{j=1}^{d} \frac{(X_{i,j} - (M_X^{(l)})_{i,j})^2}{2(S_X^{(l)})_{i,j}^2}$$

$$- \frac{1}{L} \sum_{l=1}^{L} \sum_{i=1}^{m} \sum_{j=1}^{d} \log (S_X^{(l)})_{i,j} - c \qquad (B.6)$$

To promise the acylic of $A$, the constraint Eq. (2) is adopted and the optimization problem can be formulated as below :

$$\min_{A} \quad Loss = E_{q(Z|X)}[\log p(X|Z)] + D_{KL}[q(Z|X)||p(Z)]$$
$$s.t. \quad h(A) = Tr[(I + \gamma A \odot A)^d] - d = 0$$

## References

[1] T.-Y. Chai, "Development directions of automation science and technology," *Acta Automatica Sinica*, vol. 44, no. 11, pp. 1923–1930, 2018.

[2] W. Gui, X. Chen, C. Yang, and Y. Xie, "Knowledge automation and its industrial application," *Scientia Sinica Informationis*, vol. 46, no. 8, pp. 1016–1034, 2016.

[3] L. Yao, W. Shao, and Z. Ge, "Hierarchical quality monitoring for large-scale industrial plants with big process data," *IEEE Transactions on Neural Networks and Learning Systems*, pp. 1–12, 2019.

[4] X. Yuan, Y. Gu, Y. Wang, C. Yang, and W. Gui, "A deep supervised learning framework for data-driven soft sensor modeling of industrial processes," *IEEE Transactions on Neural Networks and Learning Systems*, vol. 31, no. 11, pp. 4737–4746, 2020.

[5] Q. Sun and Z. Ge, "A survey on deep learning for data-driven soft sensors," *IEEE Transactions on Industrial Informatics*, pp. 1–1, 2021.

[6] L. Zeng and Z. Ge, "Improved population-based incremental learning of bayesian networks with partly known structure and parallel computing," *Engineering Applications of Artificial Intelligence*, vol. 95, p. 103920, 2020.

[7] D. Heckerman, "A tutorial on learning with bayesian networks," *Innovations in Bayesian networks*, pp. 33–82, 2008.

[8] D. Koller and N. Friedman, *Probabilistic graphical models: principles and techniques*. MIT press, 2009.

[9] G. Chen and Z. Ge, "Robust bayesian networks for low-quality data modeling and process monitoring applications," *Control Engineering Practice*, vol. 97, p. 104344, 2020.

[10] J. Zheng, J. Zhu, G. Chen, Z. Song, and Z. Ge, "Dynamic bayesian network for robust latent variable modeling and fault classification," *Engineering Applications of Artificial Intelligence*, vol. 89, p. 103475, 2020.

[11] Z. Liu, Z. Ge, G. Chen, and Z. Song, "Adaptive soft sensors for quality prediction under the framework of bayesian network," *Control Engineering Practice*, vol. 72, pp. 19–28, 2018.

[12] M. Scanagatta, A. Salmerón, and F. Stella, "A survey on bayesian network structure learning from data," *Progress in Artificial Intelligence*, vol. 8, no. 4, pp. 425–439, 2019.

[13] X. Zheng, B. Aragam, P. K. Ravikumar, and E. P. Xing, "Dags with no tears: Continuous optimization for structure learning," in *Advances in Neural Information Processing Systems*, S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, Eds., vol. 31. Curran Associates, Inc., 2018. [Online]. Available: https://proceedings.neurips.cc/paper/2018/file/e347c51419ffb23ca3fd5050202f9c3d-Paper.pdf

[14] P. Spirtes, C. N. Glymour, R. Scheines, and D. Heckerman, *Causation, prediction, and search*. MIT press, 2000.

[15] G. F. Cooper and E. Herskovits, "A bayesian method for the induction of probabilistic networks from data," *Machine Learning*, vol. 9, no. 4, pp. 309–347, 1992.

[16] D. Madigan, J. York, and D. Allard, "Bayesian graphical models for discrete data," *International Statistical Review / Revue Internationale de Statistique*, vol. 63, no. 2, pp. 215–232, 1995.

[17] I. Tsamardinos, L. E. Brown, and C. F. Aliferis, "The max-min hill-climbing bayesian network structure learning algorithm," *Machine learning*, vol. 65, no. 1, pp. 31–78, 2006.

[18] J. Cussens, "Gobnilp: Learning bayesian network structure with integer programming," in *International Conference on Probabilistic Graphical Models*, vol. 138. PMLR, 2020, pp. 605–608.

[19] S. Shimizu, P. O. Hoyer, A. Hyvärinen, A. Kerminen, and M. Jordan, "A linear non-gaussian acyclic model for causal discovery," *Journal of Machine Learning Research*, vol. 7, no. 10, 2006.

[20] Y. Yu, J. Chen, T. Gao, and M. Yu, "DAG-GNN: DAG structure learning with graph neural networks," in *Proceedings of the 36th International Conference on Machine Learning*, ser. Proceedings of Machine Learning Research, K. Chaudhuri and R. Salakhutdinov, Eds., vol. 97. PMLR, 09–15 Jun 2019, pp. 7154–7163.

[21] I. Ng, S. Zhu, Z. Chen, and Z. Fang, "A graph autoencoder approach to causal structure learning," *arXiv preprint arXiv:1911.07420*, 2019.

[22] Y. Wang, V. Menkovski, H. Wang, X. Du, and M. Pechenizkiy, "Causal discovery from incomplete data: a deep learning approach," *arXiv preprint arXiv:2001.05343*, 2020.

[23] L. Bottou, F. E. Curtis, and J. Nocedal, "Optimization methods for large-scale machine learning," *SIAM Review*, vol. 60, no. 2, pp. 223–311, 2018.

[24] I. Ng, A. Ghassami, and K. Zhang, "On the role of sparsity and dag constraints for learning linear dags," in *Advances in Neural Information Processing Systems*, H. Larsell, M. Ranzato, R. Hadsell, M. F. Balcan, and H. Lin, Eds., vol. 33. Curran Associates, Inc., 2020, pp. 17943–17954. [Online]. Available: https://proceedings.neurips.cc/paper/2020/file/d04d42cdf14579cd294e5079e0745411-Paper.pdf

[25] D. P. Bertsekas, "Nonlinear programming," *Journal of the Operational Research Society*, vol. 48, no. 3, pp. 334–334, 1997.

[26] J. R. Magnus, "On the concept of matrix derivative," *Journal of Multivariate Analysis*, vol. 101, no. 9, pp. 2200–2206, 2010.

[27] K. B. Petersen and M. S. Pedersen, "The matrix cookbook," nov 2012, version 20121115. [Online]. Available: http://www2.compute.dtu.dk/pubdb/pubs/3274-full.html

[28] T. Hastie, R. Tibshirani, and M. Wainwright, *Statistical learning with sparsity: the lasso and generalizations*. CRC press, 2015.

[29] J. H. Christensen and D. F. Rudd, "Structuring design computations," *AIChE Journal*, vol. 15, no. 1, pp. 94–100, 1969.

[30] R. L. Motard and A. W. Westerberg, "Exclusive tear sets for flowsheets," *AIChE Journal*, vol. 27, no. 5, pp. 725–732, 1981.

[31] P. M. Castro and I. E. Grossmann, "Generalized disjunctive programming as a systematic modeling framework to derive scheduling formulations," *Industrial & Engineering Chemistry Research*, vol. 51, no. 16, pp. 5781–5792, 2012.

[32] J. J. Downs and E. F. Vogel, "A plant-wide industrial process control problem," *Computers & Chemical Engineering*, vol. 17, no. 3, pp. 245–255, 1993.

[33] J. Kuipers, G. Moffa, and D. Heckerman, "Addendum on the scoring of gaussian directed acyclic graphical models," *Annals of Statistics*, vol. 42, no. 4, pp. 1689–1691, 2014.

[34] G. Schwarz, "Estimating the dimension of a model," *Annals of statistics*, vol. 6, no. 2, pp. 461–464, 1978.

[35] C. Squires. (2021) graphical model learning. [Online]. Available: https://github.com/uhlerlab/graphical_model_leaing/tree/a790b25bd506f50a63930ab01f1ead2f0adfd30b