

# Reconstruction of a directed acyclic graph with intervention\*

Si Peng and Xiaotong Shen

*School of Statistics, University of Minnesota*  
313 Ford Hall, 224 Church St SE  
Minneapolis, MN 55455  
e-mail: [pengx179@umn.edu](mailto:pengx179@umn.edu); [xshen@umn.edu](mailto:xshen@umn.edu)

Wei Pan

*Division of Biostatistics, University of Minnesota*  
420 Delaware St. S.E.  
Minneapolis, MN 55455  
e-mail: [panxx014@umn.edu](mailto:panxx014@umn.edu)

**Abstract:** Identification of causal relations among variables is central to many scientific investigations, as in regulatory network analysis of gene interactions and brain network analysis of effective connectivity of causal relations between regions of interest. Statistically, causal relations are often modeled by a directed acyclic graph (DAG), and hence that reconstruction of a DAG's structure leads to the discovery of causal relations. Yet, reconstruction of a DAG's structure from observational data is impossible because a DAG Gaussian model is usually not identifiable with unequal error variances. In this article, we reconstruct a DAG's structure with the help of interventional data. Particularly, we construct a constrained likelihood to regularize intervention in addition to adjacency matrices to identify a DAG's structure, subject to an error variance constraint to further reinforce the model identifiability. Theoretically, we show that the proposed constrained likelihood leads to identifiable models, thus correct reconstruction of a DAG's structure through parameter estimation even with unequal error variances. Computationally, we design efficient algorithms for the proposed method. In simulations, we show that the proposed method enables to produce a higher accuracy of reconstruction with the help of interventional observations.

**MSC2020 subject classifications:** Primary 62-09.

**Keywords and phrases:** Causal relations, constrained likelihood, intervention, reconstruction identifiability.

Received May 2020.

## Contents

1	Introduction . . . . .	4134
2	Method . . . . .	4135

---

\*The authors thank the editor, the associate editor and anonymous referees for helpful comments and suggestions. Research supported in part by NSF grants DMS-1712564, DMS-1721216, DMS-1952539, and NIH grants 1R01GM126002, 2R01HL105397, 1R01AG065636, R01AG069895.

2.1	Intervention or covariate models and variance constraint . . . .	4136
2.2	Constrained maximum likelihood . . . . .	4136
3	Computation . . . . .	4137
3.1	Optimization subject to the variance constraint . . . . .	4138
3.2	Algorithm for solving (9) . . . . .	4138
4	Theory . . . . .	4140
5	Numerical study . . . . .	4142
5.1	Simulations . . . . .	4142
5.2	Analysis of Alzheimer’s disease dataset . . . . .	4149
5.3	Analysis of cytometry data . . . . .	4150
6	Discussion . . . . .	4155
A	Technical Details . . . . .	4155
A.1	Computation details for solving (9) . . . . .	4155
A.2	Analytic updating expressions for ADMM in (20) . . . . .	4156
A.2.1	<b>A</b> -step and <b>B</b> -step . . . . .	4156
A.2.2	<b>C</b> -step . . . . .	4156
A.2.3	<b>F</b> -step . . . . .	4156
A.2.4	$\lambda$ -step and $\xi$ -step . . . . .	4157
A.3	Computation details for estimating $D_0$ . . . . .	4158
A.3.1	<b>R</b> -step . . . . .	4158
A.4	Technical proofs . . . . .	4159
	References . . . . .	4162

## 1. Introduction

Directed acyclic graph (DAG) models are useful to describe pairwise causal relations between random variables, defined by a certain Markov property [5], with each node representing one variable and each directed edge representing the corresponding pairwise causal relation. DAG models have been widely used in gene and social networks [7, 19]. To identify causal relations, intervention observations are usually collected in addition to observational attributes [16]. The central topic this article addresses is the reconstruction of a DAG model based on interventional data and pertinent issues with respect to the effect of the intervention on the reconstruction of a DAG’s structure.

In the literature, it is generally believed that interventions may help the reconstruction of a DAG’s structure, particularly when a DAG model is not identifiable from data, that is, DAGs in a Markov equivalence class are not distinguishable based on observational data alone [16]. In biological experiments, for instance, intervention occurs in a form of randomized treatments in a clinical trial or a form of gene knockdown or knockout experiments in systems biology. In such a situation, some or all system variables are controlled, permitting direction estimation of ambiguous edges connecting to these controlled variables. Yet exactly how intervention impacts reconstruction of a DAG’s structure remains unknown. Consequently, it is practically important to design a reconstruction

method for interventional data, permitting the identification of a DAG’s structure. Most existing methods for intervention [6, 10, 8] assume known intervention, that is, affected variables of the intervention are known *a priori* before data collection; see [9] for references therein. However, assuming known intervention is impractical, as in system biology, where the effect of various chemicals intervening a system cannot be precisely known. To our knowledge, one exception is a Bayesian method of [4], which is designed for a low-dimensional problem without theoretical guarantee, due to the super-exponential complexity in the number of nodes.

In this article, we propose a novel variance constraint on interventions to fully identify the DAG structure in the framework of unknown intervention. Theoretically, we show in Theorem 1 that a DAG structure is fully identifiable under the constraint, which is otherwise only possible when all the error variances are not the same [18]. Moreover, we propose a constrained maximum likelihood to seek the most efficient interventions by sparsity pursuit, leading to an identifiable reconstruction of a DAG’s structure. Computationally, we develop an efficient algorithm to solve nonconvex minimization subject to the quadratic variance constraint based on the alternating direction method of multipliers (ADMM) [2]. In simulations, we investigate the impact of the intervention on reconstruction and compare the proposed method with its counterpart without intervention. Overall, the proposed method performs well.

This article is organized into seven sections. Section 2 introduces the proposed method and discusses the issue of identifiability due to intervention, followed by the computational development in Section 3. Section 4 establishes the consistency theorems of the proposed method. Section 5 performs some simulations to study the intervention effect, and two real datasets are analyzed. Section 6 summarizes the results. Finally, the Appendix A contains technical details and proofs.

## 2. Method

Consider a causal model consisting of  $p$  random variables  $\mathbf{Y} = (Y_1, \dots, Y_p)^\top$  described by a DAG, with each node representing one variable and directed edges encoding causal relations between any two variables, where  $^\top$  denotes the transpose. This model factorizes the joint distribution of  $\mathbf{Y}$ ,  $P(\mathbf{Y})$  into a product of conditional distributions of each variable given its parents:  $P(\mathbf{Y}) = \prod_{j=1}^p P(Y_j | \mathbf{pa}_j)$ , where  $\mathbf{pa}_j$  denotes the parent set of  $Y_j$  and is defined to be empty if  $Y_j$  has no parents. This factorization is known as the local Markov property [5].

A DAG is modeled by structural equations as

$$\mathbf{Y} = \mathbf{A}\mathbf{Y} + \boldsymbol{\epsilon}, \quad \boldsymbol{\epsilon} \sim N(\mathbf{0}, \mathbf{D}), \quad (1)$$

where  $\boldsymbol{\epsilon} = (\epsilon_1, \epsilon_2, \dots, \epsilon_p)^\top$  represents the latent or unexplained error,  $\mathbf{D} = \text{Diag}(\sigma_1^2, \dots, \sigma_p^2)$  is the error covariance matrix and  $\mathbf{A} = (A_{ij})_{p \times p}$  is an adjacency matrix that uniquely determines a DAG. Here  $A_{ij} \neq 0$  encodes an

edge from node  $j$  to node  $i$ . In (1), the inverse covariance matrix of  $\mathbf{Y}$  is  $\boldsymbol{\Omega} = (\mathbf{I} - \mathbf{A})^\top \mathbf{D}^{-1} (\mathbf{I} - \mathbf{A})$ , where  $\mathbf{I}$  is the identity matrix.

When  $\sigma_1 = \cdots = \sigma_p$ , (1) is identifiable [18]. Then a DAG's structure can be reconstructed by estimating  $\mathbf{A}$ . However, when  $\sigma_1 = \cdots = \sigma_p$  breaks down, (1) is usually not identifiable, which means that  $\mathbf{A}$  is not estimable.

### 2.1. Intervention or covariate models and variance constraint

To treat non-identifiability in an observational study, consider a model consisting of  $W$  intervention variables  $\mathbf{X} = (X_1, X_2, \dots, X_W)^\top$ , where the outcome of  $\mathbf{Y}$  is observed with intervention variables  $\mathbf{X}$  that may be non-informative.

After incorporating the intervention variables into (1), we obtain that

$$\mathbf{Y} = \mathbf{A}\mathbf{Y} + \mathbf{B}\mathbf{X} + \boldsymbol{\epsilon}, \quad \boldsymbol{\epsilon} \sim N(\mathbf{0}, \mathbf{D}), \quad (2)$$

where  $\mathbf{A}$ ,  $\boldsymbol{\epsilon}$  and  $\mathbf{D}$  are defined as in (1),  $\mathbf{B} = (B_{jw})_{p \times W}$  is an unknown intervention coefficient matrix, whose  $j$ th entry  $B_{jw}$  indicates the directional strength of the intervention of  $X_w$  on  $Y_j$ . When  $B_{jw} = 0$ ;  $j = 1, \dots, p$ , there is no intervention of  $X_w$  on  $Y_i$ , and thus  $X_w$  is non-informative. Note that (2) becomes a causal model with covariates  $\mathbf{X}$ .

In the situation of unequal error variances, with the help of the intervention, we may impose constraints to achieve model identifiability, which otherwise is impossible [18]. Assume that  $\mathbf{X} \sim N(\mathbf{0}, \boldsymbol{\Sigma}_X)$ , which is independent of  $\boldsymbol{\epsilon}$ . Without loss of generality, assume that  $\boldsymbol{\Sigma}_X = \mathbf{I}$  subsequently because we can reparametrize  $\mathbf{X}$  as  $\boldsymbol{\Sigma}_X^{-1/2} \mathbf{X}$ . Then under (2),  $\mathbf{Y} \sim N(\mathbf{0}, \boldsymbol{\Omega}^{-1})$ , where  $\boldsymbol{\Omega} = (\mathbf{I} - \mathbf{A})^\top (\mathbf{B}\mathbf{B}^\top + \mathbf{D})^{-1} (\mathbf{I} - \mathbf{A})$ . In (2), we impose the variance constraint:

$$\text{Var}(\mathbf{B}\mathbf{X} + \boldsymbol{\epsilon}) = \theta \mathbf{I}, \quad \text{or } \mathbf{B}\mathbf{B}^\top + \mathbf{D} = \theta \mathbf{I}, \quad (3)$$

where  $\theta > 0$  is a parameter to be estimated.

More details are deferred to Sections 3 and 4.

**Theorem 1** (Identifiability). *Assume that  $\mathbf{X} \sim N(\mathbf{0}, \boldsymbol{\Sigma}_X)$  is independent of  $\boldsymbol{\epsilon}$  in (2), and  $A_{jk} \neq 0$  for all  $k$  which is a parent of  $j$ ;  $j = 1, 2, \dots, p$ . Under (3),  $\mathbf{A}$  is identifiable from the distribution of  $(\mathbf{Y}, \mathbf{X})$ .*

As suggested by Theorem 1,  $\mathbf{A}$  in (2) becomes identifiable when (3) is imposed on  $\mathbf{B}$ , which is otherwise impossible. Note, however, that interventions  $\mathbf{B}$  leading to identifiable  $\mathbf{A}$  may not be unique. In what is to follow, we impose a sparsity constraint to identify a most sparse  $\mathbf{B}$  in terms of the number of nonzero elements of  $\mathbf{B}$ .

### 2.2. Constrained maximum likelihood

This section estimates  $\mathbf{A}$  subject to the DAG requirement while seeking a most sparse  $\mathbf{B}$  subject to (3). Consequently, the smallest set of informative intervention variables can be identified through  $\mathbf{B}$ .

Under (2), two data matrices  $(y_{ij})_{n \times p}$  and  $(x_{iw})_{n \times W}$  are observed, with  $n$  representing the sample size. Then the negative log likelihood is

$$l(\mathbf{A}, \mathbf{B}, \mathbf{D}) = \sum_{j=1}^p \left[ -\frac{n}{2} \log \sigma_j^2 + \frac{1}{2\sigma_j^2} \sum_{i=1}^n \left( y_{ij} - \sum_{k \neq j} A_{jk} y_{ik} - \sum_{w=1}^W B_{jw} x_{iw} \right)^2 \right]. \quad (4)$$

To identify nonzero entries of  $\mathbf{A}$  and  $\mathbf{B}$ , we impose sparsity constraints to regularize:

$$\sum_{1 \leq j \neq l \leq p} I(A_{jl} \neq 0) \leq K_1, \quad \sum_{1 \leq j \leq p, 1 \leq l \leq W} I(B_{jl} \neq 0) \leq K_2, \quad (5)$$

where  $K_1$  and  $K_2$  are nonnegative integer-valued tuning parameters. Note that the constraint on  $\mathbf{B}$  removes zero entries thus zero-columns of  $\mathbf{B}$ , which can be regarded as selection of intervention variables. To reinforce the DAG requirement, we impose acyclic constraints [27] to reinforce the DAG requirement to ensure no loops to occur:

$$\lambda_{ik} + I(j \neq k) - \lambda_{jk} \geq I(A_{ij} \neq 0); i, j, k = 1, \dots, p, i \neq j, \quad (6)$$

where  $\boldsymbol{\lambda} = \{\lambda_{jl}\}_{p \times p}$  is a dual variable matrix.

For computation, we replace the indicator functions in (5) and (6) by its computational surrogate  $J_\tau(z) = \min(\frac{|z|}{\tau}, 1)$  [20] to circumvent the difficulty of non-discontinuity in optimization. This yields

$$\sum_{1 \leq j < l \leq p} J_\tau(A_{jl}) \leq K_1, \quad \sum_{1 \leq j \leq p, 1 \leq l \leq W} J_\tau(B_{jl}) \leq K_2, \quad (7)$$

$$\lambda_{ik} + I(j \neq k) - \lambda_{jk} \geq J_\tau(A_{ij}); i, j, k = 1, \dots, p, i \neq j, \quad (8)$$

where  $J_\tau(z)$  approximates the indicator function as  $\tau \rightarrow 0^+$ .

Minimizing (4) in  $(\mathbf{A}, \mathbf{B}, \mathbf{D})$  subject to (7), (8), and (3) yields the constrained maximum likelihood estimate (CMLE):

$$\begin{aligned} \min_{\mathbf{A}, \mathbf{B}, \mathbf{D}} l(\mathbf{A}, \mathbf{B}, \mathbf{D}) = \\ \sum_{j=1}^p \left[ -\frac{n}{2} \log \sigma_j^2 + \frac{1}{2\sigma_j^2} \sum_{i=1}^n \left( y_{ij} - \sum_{k \neq j} A_{jk} y_{ik} - \sum_{w=1}^W B_{jw} x_{iw} \right)^2 \right], \\ \text{subj to } \sum_{1 \leq j < l \leq p} J_\tau(A_{jl}) \leq K_1, \quad \sum_{1 \leq j \leq p, 1 \leq l \leq W} J_\tau(B_{jl}) \leq K_2, \\ \lambda_{ik} + I(j \neq k) - \lambda_{jk} \geq J_\tau(A_{ij}); i, j, k = 1, \dots, p, i \neq j, \\ \mathbf{B}\mathbf{B}^\top + \mathbf{D} = \boldsymbol{\theta}\mathbf{I}, \end{aligned} \quad (9)$$

where  $(K_1, K_2, \tau)$  are tuning parameters.

### 3. Computation

This section develops a computational strategy to solve (9). First,  $\boldsymbol{\theta}$  is estimated by  $\hat{\boldsymbol{\theta}}$  through an estimate  $\hat{\mathbf{A}}$  of  $\mathbf{A}$  from the method in [27], that is,

$$\hat{\boldsymbol{\theta}} = \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^p \left( y_{ij} - \sum_{k \neq j, k=1}^p \hat{A}_{jk} y_{ik} \right)^2, \quad (10)$$

where  $\{\hat{A}_{jk}\}_{j,k=1,\dots,p}$  are obtained by a constrained maximum likelihood estimate with the sparsity constraint and acyclic constraint, based on the structural equation model (1) with observational data  $\mathbf{Y}$  alone. Then we solve (9) with  $\theta$  replaced by  $\hat{\theta}$  using a blockwise coordinate descent alternating between two blocks  $(\mathbf{A}, \mathbf{B})$  and  $\mathbf{D}$  until convergence. In particular, the  $(\mathbf{A}, \mathbf{B})$ -block is solved via a difference convex (DC) programming followed by an alternating direction method of multipliers (ADMM), while the  $\mathbf{D}$ -block is updated by gradient descent. More details are further discussed subsequently.

### 3.1. Optimization subject to the variance constraint

To deal with the variance constraint in (9), we first consider a general constrained minimization subject to the variance constraint as follows:

$$\min_{\mathbf{B}} f(\mathbf{B}), \quad \text{subj to } \mathbf{B}\mathbf{B}^\top = \mathbf{\Lambda}, \quad (11)$$

where  $f(\mathbf{B})$  is a cost function and  $\mathbf{\Lambda}$  is a diagonal matrix.

For (11), we work with its equivalent form by introducing a dual matrix  $\mathbf{C}$  to decouple  $\mathbf{B}$  and the constraint:  $\min_{(\mathbf{B}, \mathbf{C})} f(\mathbf{B})$ , subject to  $\mathbf{C}\mathbf{C}^\top = \mathbf{\Lambda}$  and  $\mathbf{C} - \mathbf{B} = \mathbf{0}$ . Then we apply the alternating direction method of multipliers (ADMM) [2] to obtain its augmented Lagrangian:  $L_\rho(\mathbf{B}, \mathbf{C}, \mathbf{y}) = f(\mathbf{B}) + \mathbf{y}^\top \text{vec}(\mathbf{C} - \mathbf{B}) + \frac{\rho}{2} \|\mathbf{C} - \mathbf{B}\|_F^2$ , where  $\mathbf{y} \in \mathbb{R}^{pW}$  is the Lagrangian multiplier for constraint  $\mathbf{C} - \mathbf{B} = \mathbf{0}$ , and  $\rho > 0$  is the augmented Lagrangian parameter. Then it is further simplified by introducing a dual variable matrix  $\mathbf{V} = \{V_{jl}\}_{p \times W}$  to incorporate  $\mathbf{y}^\top \text{vec}(\mathbf{C} - \mathbf{B})$  into the quadratic form

$$\min_{(\mathbf{B}, \mathbf{C}, \mathbf{V})} L_\rho(\mathbf{B}, \mathbf{C}, \mathbf{V}) = f(\mathbf{B}) + \frac{\rho}{2} \|\mathbf{C} - \mathbf{B} + \mathbf{V}\|_F^2, \quad \text{subj to } \mathbf{C}\mathbf{C}^\top = \mathbf{\Lambda}. \quad (12)$$

Now we apply ADMM to iterate through three steps to solve (12) until convergence. In the  $k$ th iteration,

$$\mathbf{C}^{(k+1)} = \underset{\mathbf{C}}{\text{argmin}} \frac{\rho}{2} \|\mathbf{C} - \mathbf{B}^{(k)} + \mathbf{V}^{(k)}\|_F^2, \quad \text{subj to } \mathbf{C}\mathbf{C}^\top = \mathbf{\Lambda}, \quad (13)$$

$$\mathbf{B}^{(k+1)} = \underset{\mathbf{B}}{\text{argmin}} f(\mathbf{B}) + \frac{\rho}{2} \|\mathbf{C}^{(k+1)} - \mathbf{B} + \mathbf{V}^{(k)}\|_F^2, \quad (14)$$

$$\mathbf{V}^{(k+1)} = \mathbf{V}^{(k)} + \mathbf{C}^{(k+1)} - \mathbf{B}^{(k+1)}. \quad (15)$$

In (13)–(15), the variance constraint  $\mathbf{C}\mathbf{C}^\top = \mathbf{\Lambda}$  enters only in (13). Next we provide a closed-form solution of (13) in Lemma 1.

**Lemma 1.** *The solution of (13) can be written as  $\mathbf{C}^{(k+1)} = \mathbf{\Lambda}^{1/2} \mathbf{P}\mathbf{O}^\top$ , where a singular value decomposition of  $(\mathbf{B}^{(k)} - \mathbf{V}^{(k)})^\top \mathbf{\Lambda}^{1/2}$  gives  $\mathbf{O}\mathbf{E}\mathbf{P}^\top$ ,  $\mathbf{O} \in \mathbb{R}^{W \times p}$ ,  $\mathbf{E}, \mathbf{P} \in \mathbb{R}^{p \times p}$  are the left, diagonal, and right matrices in the decomposition.*

### 3.2. Algorithm for solving (9)

After plugging  $\hat{\theta}$  into (9), we begin with the update of the  $(\mathbf{A}, \mathbf{B})$ -block by fixing  $\mathbf{D}$  at an initial value  $\mathbf{D}_0$  and optimize (9) with regard to  $(\mathbf{A}, \mathbf{B})$ . A good

estimate of  $\mathbf{D}_0$  can be obtained by solving (9) without the variance constraint, details are given in the Appendix.

To solve (9) with a fixed  $\mathbf{D}$ , we follow [27] to convert (9) to its equivalent dual form. The procedure contains two steps. First, we apply a DC programming method and decompose the nonconvex constraint function of the nonconvex constraints (7) and (8) into a difference of two convex functions, based on which we construct a sequence of convex approximation of nonconvex constrained sets iteratively, the details are given in the Appendix. Then at the  $m$ th iteration, we solve a relaxed subproblem (16). The iteration process continues until a termination criterion is met.

The  $m$ th subproblem amounts to

$$\begin{aligned} & \min_{(\mathbf{A}, \mathbf{B}, \boldsymbol{\lambda}, \boldsymbol{\xi})} l(\mathbf{A}, \mathbf{B}) \\ & + \frac{\mu_1}{\tau} \sum_{1 \leq j \neq l \leq p} |A_{jl}| w_{jl}^{(m-1)} + \frac{\mu_2}{\tau} \sum_{1 \leq j \leq p, 1 \leq l \leq W} |B_{jl}| v_{jl}^{(m-1)}, \\ \text{subj to } & \lambda_{js} + \tau I(l \neq s) - \lambda_{ls} = |A_{jl}|_1 w_{jl}^{(m-1)} + \tau(1 - w_{jl}^{(m-1)}) + \xi_{jls}; \\ & j, l, s = 1, \dots, p, j \neq l, \xi_{jls} \geq 0, \\ & \mathbf{B}\mathbf{B}^\top + \mathbf{D} = \hat{\boldsymbol{\theta}}\mathbf{I}, \end{aligned} \quad (16)$$

where  $\boldsymbol{\xi} = \{\xi_{jls}\}_{p \times p \times p}$ ,  $\xi_{jls} \geq 0$  is a slack variable tensor, and  $w_{jl}^{(m-1)} = I(|\hat{A}_{jl}^{(m-1)}| \leq \tau)$  and  $v_{jl}^{(m-1)} = I(|\hat{B}_{jl}^{(m-1)}| \leq \tau)$  are obtained from the  $(m-1)$ th iteration.

For (16), we apply ADMM method by decoupling  $(\mathbf{A}, \mathbf{B})$  in the likelihood from the rest part of the cost function and the acyclic constraint. As in Section 3.1, we introduce dual variable tensor  $\mathbf{y} = \{y_{jls}\}_{p \times p \times p}$ , dual variable matrix  $\mathbf{U} = \{U_{jl}\}_{p \times (p+W)}$  and  $\mathbf{V} = \{V_{jl}\}_{p \times W}$ . Then we minimize the augmented Lagrangian under the variance constraint:

$$\begin{aligned} & \min_{(\mathbf{A}, \mathbf{C}, \mathbf{B}, \mathbf{F}, \boldsymbol{\lambda}, \boldsymbol{\xi}, \mathbf{y}, \mathbf{U}, \mathbf{V})} L_\rho(\mathbf{A}, \mathbf{C}, \mathbf{B}, \mathbf{F}, \boldsymbol{\lambda}, \boldsymbol{\xi}, \mathbf{y}, \mathbf{U}, \mathbf{V}) = l(\mathbf{A}, \mathbf{B}) \\ & + \frac{\mu_1}{\tau} \sum_{1 \leq j \neq l \leq p} |F_{jl}| w_{jl}^{(m-1)} + \frac{\mu_2}{\tau} \sum_{1 \leq j \leq p, 1 \leq l \leq W} |B_{jl}| v_{jl}^{(m-1)} \\ & + \sum_{1 \leq s \leq p} \sum_{1 \leq j \neq l \leq p} \frac{\rho}{2} \left( |F_{jl}| w_{jl}^{(m-1)} + \tau(1 - w_{jl}^{(m-1)}) + \xi_{jls} - \lambda_{jl} \right. \\ & \quad \left. - \tau I(l \neq s) + \lambda_{ls} + y_{jls} \right)^2 \\ & + \frac{\rho}{2} \sum_{1 \leq j, l \leq p} (A_{jl} - F_{jl} + U_{jl})^2 \\ & + \frac{\rho}{2} \sum_{1 \leq j \leq p, 1 \leq l \leq W} (B_{jl} - F_{j, l+p} + U_{j, l+p})^2 \\ & + \frac{\rho}{2} \sum_{1 \leq j \leq p, 1 \leq l \leq W} (C_{jl} - B_{jl} + V_{jl})^2, \\ & \text{subj to } \mathbf{C}\mathbf{C}^\top + \mathbf{D} = \hat{\boldsymbol{\theta}}\mathbf{I}. \end{aligned} \quad (17)$$

Again, we solve (17) over blocks  $(\mathbf{A}, \mathbf{C}, \mathbf{B}, \mathbf{F}, \boldsymbol{\lambda}, \boldsymbol{\xi}, \mathbf{y}, \mathbf{U}, \mathbf{V})$  iteratively until convergence, where analytic updating formulas are given in the Appendix. After ADMM iterations converge, we continue the DC iterations until converge, then the current  $(\mathbf{A}, \mathbf{B})$ -block is updated.

When updating the  $\mathbf{D}$ -block, we conduct a gradient descent update for each of  $(\sigma_1^2, \dots, \sigma_p^2)$  based on the current values of  $(\hat{\mathbf{A}}, \hat{\mathbf{B}})$ . The gradients for the  $j$ th dimension is  $l'_j = -\frac{n}{2\sigma_j^2} - \frac{1}{2\sigma_j^4} \sum_{i=1}^n (y_{ij} - \sum_{k \neq j} \hat{A}_{jk} y_{ik} - \sum_{w=1}^W \hat{B}_{jw} x_{iw})^2$ . Then  $\sigma_j^2$  is updated by

$$\sigma_j^2 = \sigma_j^2 - \alpha l'_j, \quad (18)$$

where  $\alpha > 0$  is the step size.

The computational strategy is summarized in Algorithm 1.

---

**Algorithm 1:** Constrained maximum likelihood

---

- Step 1. Obtain an estimate  $\hat{\theta}$  of  $\theta$  by (10), then plug  $\hat{\theta}$  into (9).  
 Step 2. Fix  $\mathbf{D}$  at an initial value  $\mathbf{D}_0$ . Set pre-specified error tolerance  $\epsilon_0$  and the maximum number of iterations  $M_0$  for blockwise coordinate descent.  
 Step 3.  $(\mathbf{A}, \mathbf{B})$ -block Initialize  $\mathbf{A}$  and  $\mathbf{B}$ . Set pre-specified error tolerance  $\epsilon_1$  and the maximum number of DC iterations  $M_1$ .  
   Step 3.1. At the  $m$ th DC iteration, compute  $\hat{\mathbf{A}}^{(m)}$  and  $\hat{\mathbf{B}}^{(m)}$  by cycling through the ADMM updating steps until convergence.  
   Step 3.2. When  $|l(\hat{\mathbf{A}}^{(m+1)}, \hat{\mathbf{B}}^{(m+1)}, \hat{\mathbf{D}}) - l(\hat{\mathbf{A}}^{(m)}, \hat{\mathbf{B}}^{(m)}, \hat{\mathbf{D}})| \leq \epsilon_1$  or  $m = M_1$ , stop the DC loop and output  $(\hat{\mathbf{A}}, \hat{\mathbf{B}}) = (\hat{\mathbf{A}}^{(m)}, \hat{\mathbf{B}}^{(m)})$ .  
 Step 4.  $(\mathbf{D})$ -block Update  $\mathbf{D}$  according to (18).  
 Step 5. Iterate through steps 3 through 4 until  $\sum_{j=1}^p |l'_j| \leq \epsilon_0$  or the number of iterations equals to  $M_0$ , stop and output  $(\hat{\mathbf{A}}, \hat{\mathbf{B}})$ .
- 

The computation complexity of Algorithm 1 is of order  $O(M_0 M_1 M_2 p^2 (p + W) + n(p + W)^2 + p(p + W)^3)$ , where  $M_0$ ,  $M_1$  and  $M_2$  are the maximum number of iterations for  $(\mathbf{A}, \mathbf{B})$ -step, DC and ADMM, respectively. For each ADMM iteration, the computation complexity is  $O(p^2(p + W))$ . The preparation phase has  $O(n(p + W)^2 + p(p + W)^3)$  complexity. In practice, the DC loop usually converges in a few iterations, which has finite termination property, c.f., Lemma 2 in [21].

#### 4. Theory

This section develops a theory quantifying the reconstruction error of the CMLE defined in (9). In particular, we first show the estimated  $\hat{\theta}$  recovers the optimal parameter estimation of the oracle estimator  $\hat{\theta}_O$ , which is defined as the maximum likelihood estimate in [27] provided that the true non-zero pattern of the DAG is given. Then we establish reconstruction consistency of the true DAG's structure defined by adjacency matrix  $\mathbf{A}$  in the presence of intervention variables  $\mathbf{X}$  in (2).

Let  $\mathbf{A}$ ,  $\mathbf{B}$ , and  $\mathbf{D}$  represent model parameters under (2). Let  $E = \{(i, j) : A_{ij} \neq 0\}$  be the set of non-zero edges in the graph  $\mathcal{G}$ , and  $|E|$  denote the size of the set. Let  $\mathbf{\Omega} = (\mathbf{I} - \mathbf{A})^\top (\mathbf{B}\mathbf{B}^\top + \mathbf{D})^{-1} (\mathbf{I} - \mathbf{A}) = (\mathbf{I} - \mathbf{A})^\top (\mathbf{I} - \mathbf{A}) / \theta$  be the precision matrix of  $\mathbf{Y}$ . In what follows, we will assume the variance constraint (3), and by Theorem 1,  $\mathbf{A}$  is identifiable from the distribution of  $(\mathbf{Y}, \mathbf{X})$ , and

thus the graph structure of  $\mathcal{G}$  is identifiable. Let  $G^0, \mathbf{A}^0, \mathbf{B}^0, \mathbf{R}^0, E^0, \mathbf{\Omega}^0, \theta^0$  denote the truth. Let  $\hat{\mathbf{A}}_O$  and  $\hat{\theta}_O$  denote the oracle estimator, or the maximum likelihood estimate provided that the true set  $E^0$  of non-zero edges is given. Let  $\hat{\mathbf{\Omega}}_O = (\mathbf{I} - \hat{\mathbf{A}}_O)^\top (\mathbf{I} - \hat{\mathbf{A}}_O) / \hat{\theta}_O$  be the oracle estimator for  $\mathbf{\Omega}$ .

For the observational model, let  $\mathbf{A}^{obs}$  be the model parameter, then the precision matrix is  $\mathbf{\Omega}^{obs} = (\mathbf{I} - \mathbf{A}^{obs})^\top (\mathbf{I} - \mathbf{A}^{obs}) / \theta$ , whose oracle estimator is  $\hat{\mathbf{\Omega}}_O^{obs} = (\mathbf{I} - \hat{\mathbf{A}}_O^{obs})^\top (\mathbf{I} - \hat{\mathbf{A}}_O^{obs}) / \hat{\theta}_O$  where  $\hat{\mathbf{A}}_O^{obs}$  is the oracle estimator of  $\mathbf{A}^{obs}$ .

The degree of reconstructability is defined as

$$C_{\min}(\mathbf{\Omega}^0) = \inf_{\{\mathbf{A} \neq \mathbf{A}^0, |E| \leq |E^0|, \mathbf{A} \text{ induces a DAG}\}} \frac{-\log(1 - h^2(\mathbf{\Omega}, \mathbf{\Omega}^0))}{\max(|E^0 \setminus E|, 1)},$$

where  $h^2(\mathbf{\Omega}, \mathbf{\Omega}^0)$  is the Hellinger distance between  $\mathbf{\Omega}$  and  $\mathbf{\Omega}^0$  under (2) and the variance constraint (3), and  $E_1 \setminus E_2$  denote the set difference between  $E_1$  and  $E_2$ . The degree of reconstructability measures the difficulty of reconstructing the graph, and we require it to be larger than a certain level in order for our proposed method to be consistent in reconstructing the graph structure. For a detailed discussion about the degree of reconstructability of a graph, c.f., [27].

**Assumption A.1** (Boundedness). For some positive constants  $M_1$  and  $M_2$ ,  $\inf_{\mathbf{\Omega}} c_{\min}(\mathbf{\Omega}) \geq M_1$ ,  $\sup_{1 \leq k \leq p} |\Omega_{kk}| \leq M_2$ , where  $c_{\min}(\mathbf{\Omega})$  is the smallest eigenvalue of  $\mathbf{\Omega}$  and  $\Omega_{kk}$  is the  $k$ th diagonal element of  $\mathbf{\Omega}$ .

**Assumption A.2** (Boundedness). For some positive constants  $M_3$  and  $M_4$ ,  $\inf_{\mathbf{\Gamma}} c_{\min}(\mathbf{\Gamma}) \geq M_3$  and  $\sup_{1 \leq k \leq p+W} |\Gamma_{kk}| \leq M_4$ , where  $\mathbf{\Pi}$  is the covariance matrix of the joint distribution of  $(\mathbf{Y}, \mathbf{X})$  and  $\mathbf{\Gamma} = \mathbf{\Pi}^{-1}$ .

**Assumption B** (Degree of reconstructability).  $C_{\min}(\mathbf{\Omega}^0) \geq 4d_0^{-1}n^{-1} \times \max(\log p, |E^0|)$ , for some positive constant  $d_0 > 0$ , say  $d_0 = \frac{2}{27} \frac{1}{963}$ .

**Assumption C.** For some positive constants  $d_1, d_2$  and  $d_3$ ,

$$h^2(\mathbf{\Omega}, \mathbf{\Omega}^0) \geq d_1 h^2(\mathbf{\Omega}_\tau, \mathbf{\Omega}^0) - d_3 p \tau^{d_2},$$

where  $\mathbf{\Omega}_\tau = (\mathbf{I} - \mathbf{A}_\tau)^\top (\mathbf{I} - \mathbf{A}_\tau) / \theta$  and  $\mathbf{A}_\tau$  is a truncated version of  $\mathbf{A}$  with its  $ij$ th element  $A_{ij} I(|A_{ij}| \geq \tau)$ .

Assumptions A.1 and A.2 concern the smallest eigenvalues and the maximum diagonal element of  $\mathbf{\Omega}$ . Under Assumption A.1, the likelihood function becomes bounded. Assumption B is a key condition for the consistency of reconstructed graph structure, we require the degree of reconstructability to be no less than a lower bound, which is related to the size of  $p$  or  $|E^0|$ . Assumption C requires the Hellinger distance to be smooth so that we approximate  $L_0$  function with its computational surrogate, the TLP function [21] to the desired level by tuning  $\tau$ .

First, we show that  $\theta$  is uniquely defined regardless of the value of  $\mathbf{B}$ .

**Lemma 2.** Under (2),  $\theta$  uniquely satisfies (3).

Then we show the optimal parameter estimation of  $\theta$  achieved by utilizing the observational data in the Theorem 2.

**Theorem 2** (Optimal parameter estimation). *Under Assumptions A.1 and C, if  $K_1 = |E^0|$  and  $\tau \leq C_{\min}(\mathbf{\Omega}^0)M_1/4p$ , then there exists a constant  $c_2 > 0$ , say  $c_2 = \frac{4}{27} \frac{1}{1926}$ , such that for any  $(n, |E^0|, p)$ ,*

$$P(\hat{\theta} \neq \hat{\theta}_O) \leq P(\hat{\mathbf{\Omega}}^{obs} \neq \hat{\mathbf{\Omega}}_O^{obs}) \leq \exp(-c_2 n C_{\min}(\mathbf{\Omega}^0) + 2 \log(p(p-1) + 1) + 3).$$

*Under Assumption B,  $P(\hat{\theta} \neq \hat{\theta}_O) \rightarrow 0$  as  $n, p, |E^0| \rightarrow \infty$ .*

The next theorem gives a reconstruction error bound, under which we obtain reconstruction consistency of the CMLE as well as its optimal parameter estimation.

**Theorem 3** (Error bound and oracle properties). *Under Assumptions A.1, A.2 and C, if  $K_1 = |E^0|$ ,  $\tau \leq C_{\min}(\mathbf{\Omega}^0)M_1/4p$ , then there exists a constant  $c_2 > 0$ , say  $c_2 = \frac{4}{27} \frac{1}{1926}$ , such that for any  $(n, |E^0|, p)$ ,*

$$P(\hat{G} \neq G^0) \leq P(\hat{\mathbf{\Omega}} \neq \hat{\mathbf{\Omega}}_O) \leq \exp(-c_2 n C_{\min}(\mathbf{\Omega}^0) + 2 \log(p(p-1) + 1) + 3).$$

*Under Assumption B,  $P(\hat{G} \neq G^0) \rightarrow 0$ ,  $\frac{Eh^2(\hat{\mathbf{\Omega}}, \mathbf{\Omega}^0)}{Eh^2(\hat{\mathbf{\Omega}}_O, \mathbf{\Omega}^0)} \rightarrow 1$  as  $n, p, |E^0| \rightarrow \infty$ .*

## 5. Numerical study

### 5.1. Simulations

This section examines the performance of the proposed method and demonstrates how intervention, as well as the variance constraint, improves the reconstructability of a DAG's structure. Seven methods are compared, including the proposed method with intervention, that without the variance constraint (3), the observational method without intervention [27], the constraint-based PC algorithm [22], the score-and-search method GES [3], and two hybrid methods, Max-Min Hill-Climbing (MMHC) [24] and ARGES [14]. For PC algorithm and GES, we use R package `pcalg`, while for MMHC we use R package `bnlearn`. For ARGES, we use the ARGES-CIG version [14] by first conducting a neighborhood selection using R package `huge` and then apply the greedy search using `pcalg`. For the other three methods, we implement in R with the main algorithm written in C, which is also available in the R package `intdag` [17] <https://cran.r-project.org/web/packages/intdag/index.html>.

Several performance metrics are used to measure the accuracy of reconstruction of a graph's skeleton as well as directionality. With respect to the skeleton of a graph, we use the false discovery rate (FDR) and false negative rate (FNR), defined as  $FDR = FP/(TP + FP)$  and  $FNR = FN/(TP + FN)$ , where TP, FP, TN, and FN denote the true positives, false positives, true negatives and false negatives, respectively. These two metrics together measure the abilities to control false discoveries, as well as false negatives. With regard to directionality, we employ the Structural Hamming Distance (SHD), defined as the minimal number of operations required to transform one DAG to the other, including edge insertions, deletions or flips, c.f., [24]. Note that a smaller SHD value indicates

two DAGs are closer to each other. To compute the SHD, we use the R-package `pcalg`. All the metrics are used on the estimation of adjacency matrix  $\mathbf{A}$ , since we focus on the reconstructability of DAG.

For tuning, PC algorithm and MMHC require one tuning parameter  $\alpha$  controlling the significance level for independence tests, yet there is no practical tuning way via a separate tuning set. In this simulation, the significance level is fixed at 0.05. This choice of  $\alpha$  seems sensible as the number of estimated edges is roughly the same as the number of edges in the true graph, as shown in the simulation. For ARGES, the first stage of neighborhood selection needs one tuning parameter corresponding to the LASSO penalization, we use the functions `huge.path()` in the R package `huge` to select the tuning range based on the data, and the function `huge.fit()` to select the tuning parameter. For the observational DAG method [27],  $\tau$  is chosen from a set  $\{0.1, 0.05, 0.01\}$ , and the sparsity regularization parameter  $\mu_1$  is chosen so that the number of estimated edges roughly ranges from 0 to 100. For our methods,  $\tau$  and  $\mu_1$  are selected similarly, and  $\mu_2 = r \times \mu_1$  with the ratio  $r$  selected from  $\{1, 2, 4, 8\}$ . For each method, the optimal tuning parameters are obtained by maximizing the predicted log-likelihood (4) based on an independent tuning set of size 1000 over a set pre-specified grid points.

In simulations, we examine a sparse neighborhood graph and a sparse graph with non-sparse neighborhoods in Examples 1 and 2, respectively. A sparse neighborhood requires each node to have sparse links, but a sparse graph does not necessarily have sparse neighborhoods. To further investigate the operating characteristics of the methods, we consider two additional situations in Examples 3 and 4, in which the true graphs satisfy the variance constraint while other settings resemble Examples 1 and 2, respectively. Finally, to compare the performances of the methods in non-sparse situations, Examples 5 and 6 are added by increasing the sampling probabilities when generating the edges. More details are given in the example settings.

**Example 1** (Sparse neighborhood). A DAG with 50 nodes is generated with a random generation mechanism as described in [11]. First, the partial ordering of these nodes is randomly generated. Second, we sample edges according to a binomial distribution with probability 0.02, where the edges are assigned to a weight 0.5. The intervention matrix  $\mathbf{B}$  is a diagonal matrix with diagonal values being 0.5s, describing a situation that each node in DAG is intervened by exactly one intervention covariate. Third, we set the error variances  $(\sigma_1^2, \dots, \sigma_p^2)^\top$  to be a sequence from 1.5 to 1 with equally spaced points. Finally, we sample  $\mathbf{X}$  from a  $p$ -dimensional normal distribution  $N(0, \mathbf{I})$  and generate  $\mathbf{Y}$  is generated according to (2).

**Example 2** (Non-sparse neighborhood). This example is modified from the previous example to generate a DAG of 50 nodes with a special structure of so-called “one-control-all”, where all directed edges are connected from the first node to the other nodes with connection strength of 0.5. Evidently, the neighborhood of the first node is not sparse, but the overall graph remains sparse. The intervention matrix  $\mathbf{B}$  and the error variances are the same as in Example 1.

TABLE 1

Averaged false discovery rate (FDR), false negative rate (FNR), and Structural Hamming Distance (SHD), as well as their standard errors for four competing methods based on 10 simulation replications in Example 1. Here “Our”, “Int”, “Obs”, “PC”, “GES”, “MMHC”, “ARGES” denote the proposed method subject to the variance constraint, the proposed method without the variance constraint, the proposed method without intervention, PC method, GES method, MMHC method and ARGES method, respectively. The best performers are in bold.

$(n, p)$	Method	FDR(A)	FNR(A)	SHD
(100,50)	Our	<b>0.17(0.03)</b>	<b>0.01(0.01)</b>	<b>10.60(2.32)</b>
(100,50)	Int	0.28(0.06)	0.02(0.02)	20.60(6.04)
(100,50)	Obs	0.39(0.20)	0.23(0.07)	39.10(22.90)
(100,50)	PC	0.29(0.06)	0.20(0.07)	21.10(5.20)
(100,50)	GES	0.62(0.05)	0.24(0.06)	64.50(10.28)
(100,50)	MMHC	0.23(0.08)	0.19(0.09)	22.60(6.08)
(100,50)	ARGES	0.50(0.00)	1.00(0.01)	50.00(0.00)
(200,50)	Our	<b>0.05(0.03)</b>	<b>0.00(0.00)</b>	<b>2.50(1.84)</b>
(200,50)	Int	0.14(0.04)	<b>0.00(0.01)</b>	8.20(2.74)
(200,50)	Obs	0.37(0.16)	0.13(0.06)	29.80(13.12)
(200,50)	PC	0.29(0.05)	0.11(0.06)	19.10(3.63)
(200,50)	GES	0.51(0.05)	0.15(0.08)	43.40(5.25)
(200,50)	MMHC	0.18(0.05)	0.09(0.05)	16.30(5.14)
(200,50)	ARGES	0.49(0.04)	0.73(0.15)	48.70(2.41)
(500,50)	Our	<b>0.01(0.02)</b>	<b>0.00(0.00)</b>	<b>0.50(0.97)</b>
(500,50)	Int	0.05(0.05)	<b>0.00(0.00)</b>	2.50(2.80)
(500,50)	Obs	0.18(0.10)	<b>0.00(0.01)</b>	11.70(7.60)
(500,50)	PC	0.28(0.05)	0.04(0.05)	19.00(3.89)
(500,50)	GES	0.45(0.06)	0.13(0.06)	35.50(8.40)
(500,50)	MMHC	0.19(0.07)	0.08(0.06)	13.90(5.53)
(500,50)	ARGES	0.25(0.06)	0.10(0.07)	18.40(4.40)

**Example 3** (Sparse neighborhood with a causal model satisfying (3)). This example is modified from Example 1, in which the intervention matrix is diagonal, and the squared diagonals are generated from a sequence from 0.5 to 1 with equally spaced points so that the variance constraint is satisfied. The true  $\theta$  is 2 in (3). Other settings remain the same as Example 1.

**Example 4** (Non-sparse neighborhood with a causal model satisfying (3)). This example is modified from Example 2, in which the intervention matrix and  $\theta = 2$  are set in the same fashion as in Example 3.

**Example 5** (Non-sparse case). This example is modified from Example 1, in which the sampling probability of the binomial distribution is increased to 0.1 when generating the edges, the rest are in the same fashion as in Example 1.

**Example 6** (Non-sparse case with a causal model satisfying (3)). This example is modified from Example 3, in which the sampling probability of the binomial distribution is increased to 0.1 when generating the edges, the rest are in the same fashion as in Example 3.

As indicated in Tables 1–4, the proposed method with the variance constraint (2) performs favorably against the other methods. In Examples 1 and 3, it performs the best in all the cases in terms of all the three evaluation metrics. In

TABLE 2

Averaged false discovery rate (FDR), false negative rate (FNR), and Structural Hamming Distance (SHD), as well as their standard errors for four competing methods based on 10 simulation replications in Example 2. Here “Our”, “Int”, “Obs”, “PC”, “GES”, “MMHC”, “ARGES” denote the proposed method subject to the variance constraint, the proposed method without the variance constraint, the proposed method without intervention, PC method, GES method, MMHC method and ARGES method, respectively. The best performers are in bold. ‘N/A’ means the method did not return a result after 24 hours. The best performers are in bold.

$(n, p)$	Method	FDR(A)	FNR(A)	SHD
(100,50)	Our	<b>0.59(0.19)</b>	<b>0.05(0.05)</b>	81.40(29.54)
(100,50)	Int	0.63(0.13)	0.07(0.07)	91.90(35.95)
(100,50)	Obs	0.67(0.16)	0.21(0.09)	115.20(65.37)
(100,50)	PC	0.84(0.06)	0.90(0.04)	63.20(2.57)
(100,50)	GES	0.61(0.05)	0.17(0.06)	66.00(10.80)
(100,50)	MMHC	0.68(0.04)	0.80(0.02)	60.40(3.47)
(100,50)	ARGES	0.50(0.00)	0.96(0.04)	<b>49.00(0.00)</b>
(200,50)	Our	0.38(0.18)	0.01(0.02)	35.20(18.56)
(200,50)	Int	0.37(0.24)	<b>0.00(0.01)</b>	39.40(29.05)
(200,50)	Obs	0.29(0.25)	0.06(0.06)	36.60(51.85)
(200,50)	PC	0.78(0.07)	0.82(0.06)	63.10(5.11)
(200,50)	GES	0.44(0.05)	0.06(0.01)	32.80(7.15)
(200,50)	MMHC	0.60(0.05)	0.69(0.02)	57.30(4.22)
(200,50)	ARGES	<b>0.07(0.02)</b>	0.23(0.09)	<b>11.90(4.12)</b>
(500,50)	Our	<b>0.07(0.16)</b>	<b>0.00(0.01)</b>	<b>6.40(16.78)</b>
(500,50)	Int	0.29(0.10)	<b>0.00(0.00)</b>	20.90(8.27)
(500,50)	Obs	0.28(0.18)	<b>0.00(0.00)</b>	22.50(15.20)
(500,50)	PC	N/A	N/A	N/A
(500,50)	GES	0.38(0.13)	0.03(0.02)	25.90(20.59)
(500,50)	MMHC	N/A	N/A	N/A
(500,50)	ARGES	0.11(0.14)	0.04(0.01)	7.30(15.02)

Example 2, it achieves the top performances in terms of FDR and FNR when  $n = 100$ , in terms of FNR when  $n = 200$  and performs the best in terms of all the three metrics when  $n = 500$ . In Example 4, it also achieves the top performances in terms of FDR when  $n = 100$ , in terms of FDR and SHD when  $n = 200$  and in terms of FNR when  $n = 500$ . Interestingly, in Examples 2 and 4, PC algorithm and MMHC do not handle “non-sparse neighborhood” structures well as they fail to return results within 24 hours in the case of  $n = 500$ . Also, when  $n$  is small, ARGES performs poorly, almost identifying no directed edges, such as in the cases of  $n = 100$  of Examples 1–4 and  $n = 200$  of Examples 3 and 4. One explanation is that when  $n$  is small the tuning process of the neighborhood selection step tends to select a large penalty, resulting in a very sparse conditional independent graph (CIG).

Overall, both the variance constraint (3) and intervention lead to improvements of the reconstruction accuracy across all the situations. Specifically, the proposed method has an average amount of improvement (10.00, 10.50, 8.20, 6.80) in terms of SHD over the intervention method without the variance constraint when  $n = 100$  in all the four examples, respectively. The amount of improvement becomes (5.70, 4.20, 8.30, 6.70) when  $n$  increases to 200 and (2.00, 14.50, 2.90, 19.00) when  $n = 500$ . The improvement remains noticeable in Examples 2 and 4

TABLE 3

Averaged false discovery rate (FDR), false negative rate (FNR), and Structural Hamming Distance (SHD), as well as their standard errors for four competing methods based on 10 simulation replications in Example 3. Here “Our”, “Int”, “Obs”, “PC”, “GES”, “MMHC”, “ARGES” denote the proposed method subject to the variance constraint, the proposed method without the variance constraint, the proposed method without intervention, PC method, GES method, MMHC method and ARGES method, respectively. The best performers are in bold. ‘NaN’ in the FDR column means the ARGES method does not have any discoveries in all the replications.

$(n, p)$	Method	FDR(A)	FNR(A)	SHD
(100,50)	Our	<b>0.11(0.05)</b>	<b>0.04(0.04)</b>	<b>7.60(3.41)</b>
(100,50)	Int	0.22(0.08)	0.05(0.04)	15.80(5.63)
(100,50)	Obs	0.44(0.08)	0.14(0.05)	36.70(9.48)
(100,50)	PC	0.30(0.05)	0.21(0.07)	22.80(4.54)
(100,50)	GES	0.61(0.05)	0.21(0.09)	63.30(7.97)
(100,50)	MMHC	0.27(0.09)	0.22(0.11)	26.50(8.61)
(100,50)	ARGES	NaN	1.00(0.00)	50.00(0.00)
(200,50)	Our	<b>0.04(0.06)</b>	<b>0.00(0.00)</b>	<b>2.30(3.40)</b>
(200,50)	Int	0.16(0.11)	<b>0.00(0.01)</b>	10.60(9.07)
(200,50)	Obs	0.25(0.06)	0.03(0.04)	16.80(4.92)
(200,50)	PC	0.27(0.06)	0.12(0.05)	18.00(4.14)
(200,50)	GES	0.53(0.04)	0.17(0.05)	44.70(6.43)
(200,50)	MMHC	0.22(0.07)	0.08(0.06)	18.10(5.34)
(200,50)	ARGES	0.50(0.00)	1.00(0.00)	50.00(0.00)
(500,50)	Our	<b>0.01(0.01)</b>	<b>0.00(0.00)</b>	<b>0.40(0.52)</b>
(500,50)	Int	0.06(0.02)	<b>0.00(0.00)</b>	3.30(1.34)
(500,50)	Obs	0.17(0.05)	0.02(0.02)	10.00(3.40)
(500,50)	PC	0.27(0.04)	0.05(0.03)	18.10(3.31)
(500,50)	GES	0.42(0.08)	0.10(0.06)	32.00(8.51)
(500,50)	MMHC	0.14(0.04)	0.03(0.03)	8.80(3.01)
(500,50)	ARGES	0.43(0.06)	0.46(0.16)	41.80(6.07)

even when  $n$  is large. Interestingly, a DAG’s structure can be well-reconstructed even without the variance constraint, particularly when  $n$  is large, for example,  $n = 500$  in Examples 1 and 3. Moreover, the proposed method has an average improvement of (28.50, 33.80, 29.10, 40.20) in SHD over the observational method when  $n = 100$ , (27.30, 1.40, 14.50, 60.60) when  $n = 200$  and (11.20, 16.10, 9.60, 21.70) when  $n = 500$ . The amount of improvement is significant across all the cases except the case of  $n = 200$  in Example 2.

In Examples 2 and 4, with a non-sparse neighborhood structure, the estimation becomes more challenging. In such situations, we added the simulations with  $n = 1000$ , as shown in Tables 5–6, the proposed method with the variance constraint (2) outperforms all the competitors, with SHD very close to 0 in Example 2 and exactly 0 in Example 4. As a comparison, PC and MMHC fail to return results within 24 hours, GES and ARGES perform even worse than the cases with  $n = 500$ . For large sample sizes, GES and ARGES tend to have more false positives, with estimated graph structure failing to determine the directions of the edges between the first node and the rest.

As demonstrated in Tables 7–8, the benefits of the proposed methods are evident against the competitors in the non-sparse cases of Examples 5 and 6,

TABLE 4

Averaged false discovery rate (FDR), false negative rate (FNR), and Structural Hamming Distance (SHD), as well as their standard errors for four competing methods based on 10 simulation replications in Example 4. Here “Our”, “Int”, “Obs”, “PC”, “GES”, “MMHC”, “ARGES” denote the proposed method subject to the variance constraint, the proposed method without the variance constraint, the proposed method without intervention, PC method, GES method, MMHC method and ARGES method, respectively. The best performers are in bold. ‘N/A’ means the method did not return a result after 24 hours. The best performers are in bold. ‘NaN’ in the FDR column means the ARGES method does not have any discoveries in all the replications.

$(n, p)$	Method	FDR(A)	FNR(A)	SHD
(100,50)	Our	<b>0.65(0.09)</b>	0.16(0.09)	93.40(41.51)
(100,50)	Int	0.66(0.12)	<b>0.15(0.05)</b>	100.20(49.61)
(100,50)	Obs	0.76(0.02)	0.22(0.08)	133.60(20.30)
(100,50)	PC	0.87(0.05)	0.91(0.04)	65.70(3.86)
(100,50)	GES	0.63(0.04)	0.24(0.06)	68.70(9.83)
(100,50)	MMHC	0.68(0.03)	0.79(0.02)	60.20(2.70)
(100,50)	ARGES	NaN	1.00(0.00)	<b>49.00(0.00)</b>
(200,50)	Our	<b>0.27(0.15)</b>	0.20(0.03)	<b>21.30(14.74)</b>
(200,50)	Int	0.33(0.13)	<b>0.02(0.04)</b>	28.00(15.24)
(200,50)	Obs	0.61(0.09)	0.05(0.05)	81.90(29.97)
(200,50)	PC	0.82(0.09)	0.86(0.06)	63.60(5.72)
(200,50)	GES	0.46(0.04)	0.06(0.01)	34.60(7.18)
(200,50)	MMHC	0.61(0.04)	0.70(0.02)	57.30(3.71)
(200,50)	ARGES	0.50(0.00)	0.97(0.04)	49.00(0.00)
(500,50)	Our	0.14(0.05)	<b>0.00(0.00)</b>	8.10(2.88)
(500,50)	Int	0.29(0.24)	<b>0.00(0.01)</b>	27.10(23.13)
(500,50)	Obs	0.38(0.05)	<b>0.00(0.00)</b>	29.80(5.85)
(500,50)	PC	N/A	N/A	N/A
(500,50)	GES	0.35(0.06)	0.04(0.00)	17.90(4.95)
(500,50)	MMHC	N/A	N/A	N/A
(500,50)	ARGES	<b>0.05(0.02)</b>	0.05(0.02)	<b>2.80(0.92)</b>

TABLE 5

Averaged false discovery rate (FDR), false negative rate (FNR), and Structural Hamming Distance (SHD), as well as their standard errors for four competing methods based on 10 simulation replications in Example 2. Here “Our”, “Int”, “Obs”, “PC”, “GES”, “MMHC”, “ARGES” denote the proposed method subject to the variance constraint, the proposed method without the variance constraint, the proposed method without intervention, PC method, GES method, MMHC method and ARGES method, respectively. The best performers are in bold. ‘N/A’ means the method did not return a result after 24 hours. The best performers are in bold.

$(n, p)$	Method	FDR(A)	FNR(A)	SHD
(1000,50)	Our	<b>0.02(0.02)</b>	<b>0.00(0.00)</b>	<b>1.20(1.23)</b>
(1000,50)	Int	0.09(0.06)	<b>0.00(0.00)</b>	5.00(3.40)
(1000,50)	Obs	0.17(0.12)	<b>0.00(0.00)</b>	11.10(8.37)
(1000,50)	PC	N/A	N/A	N/A
(1000,50)	GES	0.58(0.01)	<b>0.00(0.00)</b>	58.80(2.20)
(1000,50)	MMHC	N/A	N/A	N/A
(1000,50)	ARGES	0.46(0.15)	0.00(0.01)	44.40(14.90)

TABLE 6

Averaged false discovery rate (FDR), false negative rate (FNR), and Structural Hamming Distance (SHD), as well as their standard errors for four competing methods based on 10 simulation replications in Example 4. Here “Our”, “Int”, “Obs”, “PC”, “GES”, “MMHC”, “ARGES” denote the proposed method subject to the variance constraint, the proposed method without the variance constraint, the proposed method without intervention, PC method, GES method, MMHC method and ARGES method, respectively. The best performers are in bold. ‘N/A’ means the method did not return a result after 24 hours. The best performers are in bold.

$(n, p)$	Method	FDR(A)	FNR(A)	SHD
(1000,50)	Our	<b>0.00(0.00)</b>	<b>0.00(0.00)</b>	<b>0.00(0.00)</b>
(1000,50)	Int	0.04(0.03)	0.00(0.01)	2.10(1.79)
(1000,50)	Obs	0.07(0.14)	0.00(0.01)	5.20(10.00)
(1000,50)	PC	N/A	N/A	N/A
(1000,50)	GES	0.57(0.01)	<b>0.00(0.00)</b>	58.00(2.36)
(1000,50)	MMHC	N/A	N/A	N/A
(1000,50)	ARGES	0.50(0.00)	<b>0.00(0.0)</b>	49.10(0.32)

TABLE 7

Averaged false discovery rate (FDR), false negative rate (FNR), and Structural Hamming Distance (SHD), as well as their standard errors for four competing methods based on 10 simulation replications in Example 5. Here “Our”, “Int”, “Obs”, “PC”, “GES”, “MMHC”, “ARGES” denote the proposed method subject to the variance constraint, the proposed method without the variance constraint, the proposed method without intervention, PC method, GES method, MMHC method and ARGES method, respectively. The best performers are in bold. ‘NaN’ in the FDR column means the ARGES method does not have any discoveries in all the replications.

$(n, p)$	Method	FDR(A)	FNR(A)	SHD
(200,50)	Our	0.34(0.11)	<b>0.02(0.21)</b>	<b>135.80(66.63)</b>
(200,50)	Int	0.47(0.14)	0.21(0.34)	211.20(74.60)
(200,50)	Obs	0.53(0.01)	0.04(0.02)	277.30(9.79)
(200,50)	PC	<b>0.27(0.07)</b>	0.78(0.03)	203.40(7.37)
(200,50)	GES	0.86(0.02)	0.63(0.06)	604.10(23.94)
(200,50)	MMHC	0.60(0.07)	0.90(0.02)	236.60(4.84)
(200,50)	ARGES	NaN	1.00(0.00)	250.00(0.00)
(500,50)	Our	<b>0.18(0.11)</b>	<b>0.01(0.02)</b>	<b>61.70(43.24)</b>
(500,50)	Int	0.36(0.28)	0.21(0.32)	175.20(198.45)
(500,50)	Obs	0.44(0.07)	0.11(0.21)	194.40(41.48)
(500,50)	PC	0.19(0.04)	0.71(0.02)	187.50(5.70)
(500,50)	GES	0.83(0.03)	0.64(0.05)	505.80(31.71)
(500,50)	MMHC	0.55(0.06)	0.87(0.02)	225.70(6.58)
(500,50)	ARGES	NaN	1.00(0.00)	250.00(0.00)
(1000,50)	Our	<b>0.02(0.01)</b>	0.03(0.04)	<b>10.60(10.51)</b>
(1000,50)	Int	0.09(0.06)	<b>0.03(0.02)</b>	29.20(19.47)
(1000,50)	Obs	0.10(0.11)	0.04(0.04)	39.30(40.73)
(1000,50)	PC	0.18(0.04)	0.69(0.02)	181.30(7.24)
(1000,50)	GES	0.80(0.02)	0.71(0.03)	394.40(15.90)
(1000,50)	MMHC	0.55(0.05)	0.85(0.02)	222.20(6.92)
(1000,50)	ARGES	NaN	1.00(0.00)	250.00(0.00)

TABLE 8

Averaged false discovery rate (FDR), false negative rate (FNR), and Structural Hamming Distance (SHD), as well as their standard errors for four competing methods based on 10 simulation replications in Example 6. Here “Our”, “Int”, “Obs”, “PC”, “GES”, “MMHC”, “ARGES” denote the proposed method subject to the variance constraint, the proposed method without the variance constraint, the proposed method without intervention, PC method, GES method, MMHC method and ARGES method, respectively. The best performers are in bold. ‘NaN’ in the FDR column means the ARGES method does not have any discoveries in all the replications.

$(n, p)$	Method	FDR(A)	FNR(A)	SHD
(200,50)	Our	0.39(0.10)	<b>0.04(0.03)</b>	<b>164.20(69.71)</b>
(200,50)	Int	0.50(0.06)	<b>0.04(0.02)</b>	243.60(50.96)
(200,50)	Obs	0.51(0.15)	0.16(0.21)	273.10(126.29)
(200,50)	PC	<b>0.26(0.08)</b>	0.78(0.02)	202.30(6.77)
(200,50)	GES	0.86(0.03)	0.63(0.06)	588.30(27.56)
(200,50)	MMHC	0.60(0.07)	0.90(0.02)	237.70(6.13)
(200,50)	ARGES	NaN	1.00(0.00)	250.00(0.00)
(500,50)	Our	<b>0.14(0.04)</b>	<b>0.00(0.01)</b>	<b>41.50(15.35)</b>
(500,50)	Int	0.27(0.16)	0.10(0.26)	104.80(83.12)
(500,50)	Obs	0.36(0.10)	0.02(0.02)	148.80(56.29)
(500,50)	PC	0.16(0.03)	0.69(0.02)	178.80(5.25)
(500,50)	GES	0.83(0.02)	0.65(0.04)	485.20(26.17)
(500,50)	MMHC	0.59(0.07)	0.88(0.02)	231.30(5.42)
(500,50)	ARGES	NaN	1.00(0.00)	250.00(0.00)
(1000,50)	Our	<b>0.02(0.01)</b>	<b>0.02(0.03)</b>	<b>8.60(6.50)</b>
(1000,50)	Int	0.10(0.04)	0.04(0.02)	31.00(14.91)
(1000,50)	Obs	0.13(0.06)	0.04(0.03)	47.40(20.37)
(1000,50)	PC	0.18(0.05)	0.68(0.02)	180.00(7.39)
(1000,50)	GES	0.80(0.04)	0.71(0.06)	382.80(23.86)
(1000,50)	MMHC	0.57(0.07)	0.85(0.03)	223.50(7.15)
(1000,50)	ARGES	NaN	1.00(0.00)	250.00(0.00)

where each node in the graph has 10 edges instead of 2 edges on average. The proposed method with the variance constraint continues to perform the best in all cases in the three metrics, except the case of  $n = 200$  in which PC performs better in FDR. When the sample size increases from 200 to 500 and then to 1000, our methods exhibit a larger amount of improvement, as the SHD reduces from 135.80 to 61.70 and to 10.60 in Example 5 while from 164.20 to 41.50 and to 8.60 in Example 6. By comparison, ARGES fails to produce a meaningful graph in all cases, GES yields very large values of FDR, FNR, and SHD. Whereas PC and MMHC do perform better than GES, their performances do not improve much at all when the sample size increases.

In conclusion, intervention clearly has a positive impact on the reconstructability of a DAG through interventional covariates. With the variance constraint, the accuracy of reconstruction can be further enhanced.

## 5.2. Analysis of Alzheimer’s disease dataset

This section applies our proposed method to analyze the Alzheimer’s disease dataset [25], where 8560 gene expressions were collected for 176 Alzheimer’s dis-

ease patients and 187 healthy participants. Our primary goal is to reconstruct a causal network of Alzheimer’s disease-related genes with the help of intervention, and compare the DAG structures for the patient and control groups, for identifying regulatory gene-gene interactions that differentiate these two patient groups.

Biologically, transcription factors (TFs) are proteins controlling the transcription process. By regulating genes, TFs ensure that target genes have right expressions in a cell and during a biological organism. In fact, a TF binds to its target DNA sequence and thus can be mapped to specific genes. For our purpose, we use the database [23] to extract a list of human DNA binding TFs and map them to genes in the dataset, resulting in 1031 mapped TF genes. Then we treat TF genes as intervention covariates to facilitate reconstruction of causal relations encoded by the gene networks.

To identify TF gene expressions associated with Alzheimer’s disease in our dataset, we examine the KEGG database [12]. There, 168 genes in the Alzheimer’s disease pathway, among which the expressions of 99 genes are mapped to the pathway. Among the mapped genes, we perform a two-sample t-test to obtain the significant ones between the patient and control groups, resulting in 43 selected genes at the significance level 0.05.

After pre-processing, we apply the proposed method to both the groups to reconstruct two DAG networks for the disease and control groups with 176 and 187 subjects, involving  $p = 43$  genes as causal variables and  $W = 1031$  TFs as intervention variables, where the tuning parameters of the method are estimated by a five-fold cross validation. As shown in Figure 1, there are 29 and 33 estimated directed connections in the patient and control groups, respectively, with 11 shared common directed connections. Moreover, the two networks share some common sub-structures, particularly, we can find common directed connections from NDUFAB1 to ATP5A1, ATP5G3, and COX7A2. However, different sub-structures are also revealed. In the patient group, several genes, including SDHD, NDUFA1, NDUFAB1, ATP5G3, ATP5A1 and ATP5C1 have directed connections to SDHB, whereas in the control group, only two of them are connected to SDHB. This suggests that the differences in the two DAG networks reflect the disparity in the gene regulatory relations between an Alzheimer’s disease subject and a healthy subject.

Biologically, our estimated directed connections match the known biological pathway of Alzheimer’s disease in the KEGG database. For example, the estimated directed connections between NDU-genes, SDH-genes, ATP-genes and COX-genes match the pathway of the electron transport chain in mitochondria. Also, the estimated connection from CALM3 to PPP3CB matches the biological pathway from calcium-modulated protein(CaM) to protein phosphatase 3 catalytic subunit alpha(PP3CA).

### 5.3. Analysis of cytometry data

This section applies the proposed method to analyze the flow cytometry data in [19] to reconstruct causal relations between phosphorylated proteins and phos-

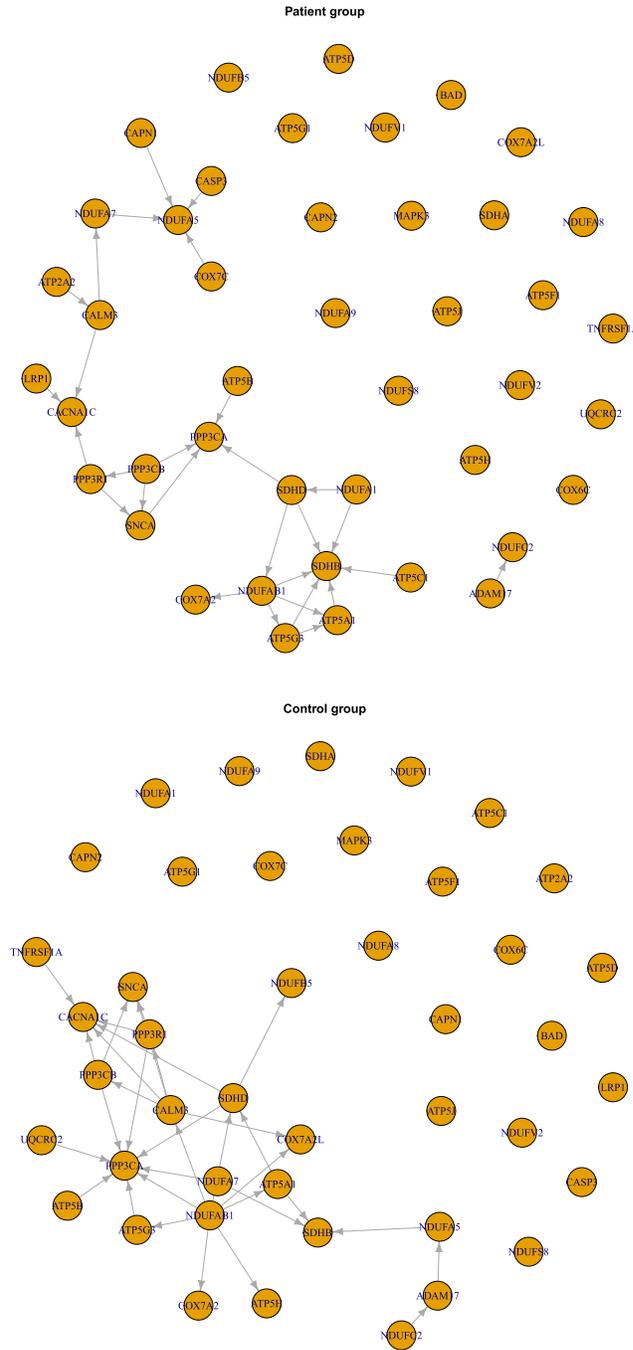


FIG 1. Top: DAG network for Alzheimer's disease patient group. There are 38 directed connections. Bottom: DAG network for healthy participant group. There are 30 directed connections, 11 of which are shared by the DAG network of patient group.

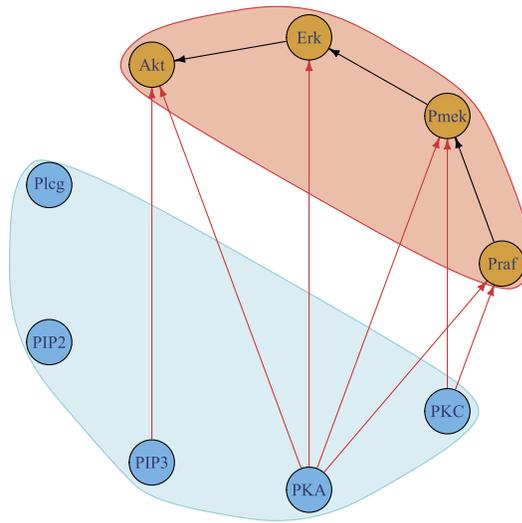


FIG 2. Consensus intracellular signaling network as benchmark. Four nodes in brown are intervened by five intervention nodes marked in blue, with directed links indicating their intervention directions by solid red lines. Note that no links are present between intervention nodes.

pholipids in human primary naive CD4+ T cells of the immune system. These cells were perturbed with molecular interventions to drive the ordering of connections in an intracellular signaling network. The original flow cytometry data contains  $n = 7466$  cell measurements, each consisting of the amount of  $p = 11$  proteins and phospholipids. The data is collected from nine different experiments in which different components in the network are intervened, either by stimulatory cues or inhibitory interventions. For our analysis, we use the continuous version of the original data [19]. Our objective is to reconstruct this network while the consensus network [19] is used as a benchmark for discovery.

For interventional data, we pre-process by selecting nine out of the eleven components, four of which form a DAG with three directed links while remaining five serve as intervention nodes; see Figure 2 for a display of an enlarge network with the nine nodes, known as consensus network, where links between the five intervention nodes are purposely removed for the intervention purpose. Note that the removal of these links does not affect our analysis, because each node is independent of its non-descendants given its direct parents by the local Markov property. Among the five intervention nodes, only three of them are informative with the other two Plcg and PIP2 having no effects on the four nodes to be intervened.

We fit our proposed method with and without interventions. For tuning, we use one-tenth of the samples for training, and the rest for tuning. As displayed in Figures 3 and 4, the proposed method with intervention correctly identi-

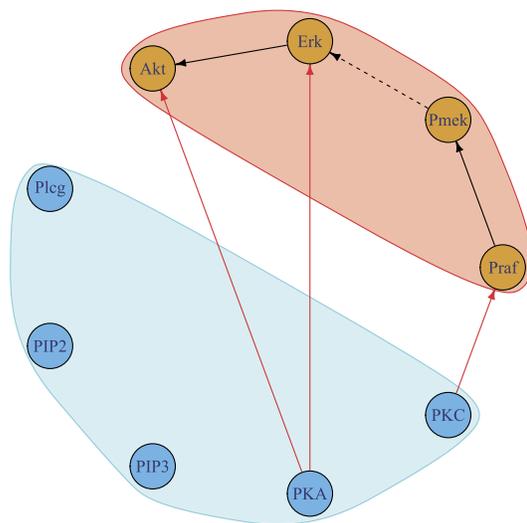


FIG 3. Intracellular signaling network reconstructed by the proposed method with intervention. Two directed links in solid black are correctly identified, denoted by solid black lines, while the directed link from *Pmek* to *Erk* is missed, denoted by a dashed black line. Three intervention links are identified by our model (denoted by solid red lines), all of which are present in the consensus network.

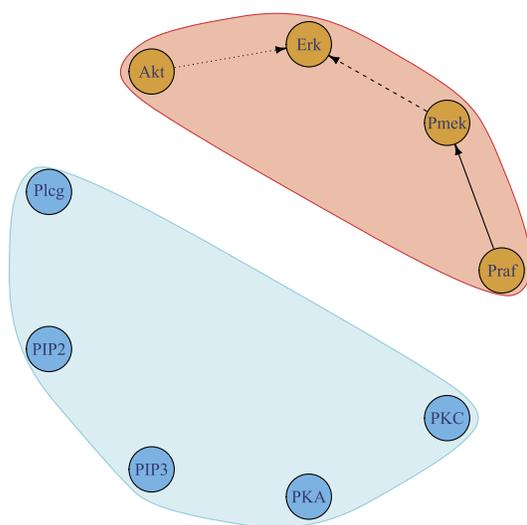


FIG 4. Intracellular signaling network estimated by the model without intervention. One directed connection is correctly identified, denoted by the solid black line, while another connection is reversed, denoted by dotted black line, and the directed connection from *Pmek* to *Erk* is missed, denoted by dashed black line.

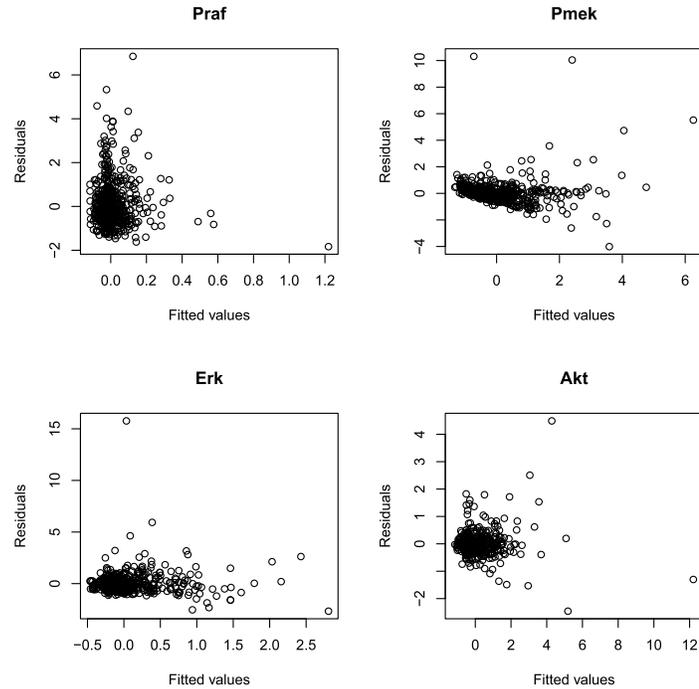


FIG 5. Residual plots from our intervention model. Each subplot corresponds to one node in the DAG network.

fies the directed link from **Erk** to **Akt**, while the observational method alters the direction oppositely. Both the methods enable to reconstruct the directed link from **Praf** to **Pmek**, whereas they both miss the link from **Pmek** to **Erk**. Consequently, the proposed method with intervention identifies more correct directed links than that without intervention. With respect to the intervention effects, our method identifies three intervention links in the network, which rule out the true non-informative intervention nodes **Plcg** and **PIP2**. **PIP3** is also non-informative in our estimation, and a possible reason is that **Akt** is already intervened by **PKA**, and thus the intervention effect of **PIP3** on **Akt** is unnecessary. Similarly, the intervention from **PKC** to **Praf** is identified by the proposed method, while both **PKA** and **PKC** have interventions on **Praf** in the consensus network. Finally, the proposed method yields a sparse intervention pattern.

One plausible explanation of missing the link from **Pmek** to **Erk** is that the linear causal model fails to capture nonlinear functional relations among cytometry measurements of proteins **Praf**, **Pmek**, **Erk**, and **Akt**, as evident from nonlinear patterns revealed by their residual plots of the structural equation model in Figure 5. Another possibility is that the five nodes on the bottom-left in Figure 3 are not originally designed as intervention nodes in the experiments.

## 6. Discussion

In this paper, a constrained maximum likelihood method is proposed to reconstruct the structure of a DAG using interventional data and efficient ADMM algorithms are developed to solve the optimization problems. In particular, a novel variance constraint is introduced and leverages the information in the interventional data to improve identifiability. Theories are established and it is shown that with the introduction of our variance constraint, the graphical structure is shown to be fully identifiable under some mild assumptions. The theoretical results are also demonstrated in our simulation results, as our proposed method performs well against competitors and can reconstruct the DAG more accurately than the observational method.

### Appendix A: Technical Details

#### A.1. Computation details for solving (9)

For (9), after plugging in  $\hat{\theta}$  and fix  $\mathbf{D}$ , we proceed in two steps. First, we relax nonconvex constraints (7) and (8) using a sequence of convex approximations involving  $2^p + 1$  linear constraints, where each approximation is refined iteratively. Then we solve each subproblem by employing a constrained alternating direction method of multipliers to estimate. The underlying process iterates until convergence.

For convex relaxation of nonconvex constraints (7) and (8) in (9), we employ difference convex (DC) programming similar to [27]. In particular, we decompose  $J_\tau$  into a difference of two convex functions:  $J_\tau(z) = S_1(z) - S_2(z) \equiv \min(\frac{|z|}{\tau}, 1) = \frac{|z|}{\tau} - \max(\frac{|z|}{\tau} - 1, 0)$ . On this ground, we construct a sequence of convex approximating sets iteratively by replacing  $S_2$  in the decomposition at iteration  $m$  by its affine majorization at iteration  $m - 1$ . (9) becomes

$$\begin{aligned}
 & \min_{(\mathbf{A}, \mathbf{B}, \boldsymbol{\lambda})} l(\mathbf{A}, \mathbf{B}, \boldsymbol{\lambda}) \\
 \text{subj to } & \frac{1}{\tau} \sum_{1 \leq j \neq l \leq p} |A_{jl}| w_{jl}^{(m-1)} \leq K_1 - \sum_{1 \leq j < l \leq p} (1 - w_{jl}^{(m-1)}), \\
 & \frac{1}{\tau} \sum_{1 \leq j \leq p, 1 \leq l \leq W} |B_{jl}| v_{jl}^{(m-1)} \leq K_2 - \sum_{1 \leq j \leq p, 1 \leq l \leq W} (1 - v_{jl}^{(k, m-1)}), \\
 & \lambda_{js} + \tau I(l \neq s) - \lambda_{ls} \geq |A_{jl}|_1 w_{jl}^{(m-1)} + \tau(1 - w_{jl}^{(m-1)}); \\
 & j, l, s = 1, \dots, p, j \neq l, \\
 & \mathbf{B}\mathbf{B}^\top + \mathbf{D} = \hat{\theta}\mathbf{I},
 \end{aligned} \tag{19}$$

where  $w_{jl}^{(m-1)} = I(\|\hat{A}_{jl}^{(m-1)}\|_1 \leq \tau)$ ,  $v_{jl}^{(m-1)} = I(|\hat{B}_{jl}^{(m-1)}| \leq \tau)$ ;  $1 \leq i, j \leq p$ , and  $(\hat{\mathbf{A}}^{(m-1)}, \hat{\mathbf{B}}^{(m-1)})$  is the solution at iteration  $m - 1$ .

Now consider a regularization version of (19), with a slack variable  $\boldsymbol{\xi}$  added to the inequality constraint, yielding (16).

### A.2. Analytic updating expressions for ADMM in (20)

At ADMM iteration step  $s + 1$ , the updating formula are

$$\begin{aligned}
\mathbf{A}^{(s+1)} &= \operatorname{argmin}_{\mathbf{A}} L_{\rho}(\mathbf{A}, \mathbf{C}^{(s)}, \mathbf{B}^{(s)}, \boldsymbol{\lambda}^{(s)}, \boldsymbol{\xi}^{(s)}, \mathbf{y}^{(s)}, \mathbf{U}^{(s)}, \mathbf{V}^{(s)}), \\
\mathbf{C}^{(s+1)} &= \operatorname{argmin}_{\mathbf{C}} L_{\rho}(\mathbf{A}^{(s+1)}, \mathbf{C}, \mathbf{B}^{(s)}, \boldsymbol{\lambda}^{(s)}, \boldsymbol{\xi}^{(s)}, \mathbf{y}^{(s)}, \mathbf{U}^{(s)}, \mathbf{V}^{(s)}) \\
&\quad \text{subj to } \mathbf{C}\mathbf{C}^{\top} + \hat{\mathbf{D}} = \sigma^2 \mathbf{I}, \\
\mathbf{B}^{(s+1)} &= \operatorname{argmin}_{\mathbf{B}} L_{\rho}(\mathbf{A}^{(s+1)}, \mathbf{C}^{(s+1)}, \mathbf{B}, \boldsymbol{\lambda}^{(s)}, \boldsymbol{\xi}^{(s)}, \mathbf{y}^{(s)}, \mathbf{U}^{(s)}, \mathbf{V}^{(s)}), \\
\boldsymbol{\lambda}^{(s+1)} &= \operatorname{argmin}_{\boldsymbol{\lambda}} L_{\rho}(\mathbf{A}^{(s+1)}, \mathbf{C}^{(s+1)}, \mathbf{B}^{(s+1)}, \boldsymbol{\lambda}, \boldsymbol{\xi}^{(s)}, \mathbf{y}^{(s)}, \mathbf{U}^{(s)}, \mathbf{V}^{(s)}), \\
\boldsymbol{\xi}^{(s+1)} &= \operatorname{argmin}_{\boldsymbol{\xi}} L_{\rho}(\mathbf{A}^{(s+1)}, \mathbf{C}^{(s+1)}, \mathbf{B}^{(s+1)}, \boldsymbol{\lambda}^{(s+1)}, \boldsymbol{\xi}, \mathbf{y}^{(s)}, \mathbf{U}^{(s)}, \mathbf{V}^{(s)}), \\
&\quad \text{subj to } \xi_{ijk} \geq 0; i, j, k = 1, \dots, p, j \neq k, \\
y_{ijk}^{(s+1)} &= y_{ijk}^{(s)} + (|F_{ij}^{(s+1)}| + \xi_{ijk}^{(s+1)} - \tau \lambda_{ik}^{(s+1)} - \tau I(j \neq k) + \tau \lambda_{jk}^{(s+1)}), \\
U_{jl}^{(s+1)} &= U^{(s)} + (A_{jl}^{(s+1)} - F_{jl}^{(s+1)}), \quad 1 \leq j, l \leq p, \\
U_{j,l+p}^{(s+1)} &= U^{(s)} + (B_{jl}^{(s+1)} - F_{j,l+p}^{(s+1)}), \quad 1 \leq j \leq p, 1 \leq l \leq W, \\
V_{jl}^{(s+1)} &= V_{jl}^{(s)} + (C_{jl}^{(s+1)} - B_{jl}^{(s+1)}), \quad 1 \leq j \leq p, 1 \leq l \leq W. \tag{20}
\end{aligned}$$

#### A.2.1. A-step and B-step

For simplicity, denote  $\mathbf{H} = (\mathbf{A}, \mathbf{B})$  be the concatenation of adjacency matrix and intervention matrix, let  $\mathbf{Z} = (\mathbf{Y}, \mathbf{X})$  be the concatenated data matrices. For each row of  $\mathbf{H}$ , the optimization problem is summarized as  $\min_{\mathbf{H}_{j,j^-}} \frac{1}{2} \|\mathbf{z}_j - \mathbf{Z}_{j^-} \mathbf{H}_{j,j^-}^{\top}\|^2 + \frac{\rho}{2} \|\mathbf{H}_{j,j^-} - \mathbf{F}_{j,j^-}^{(s)} + \mathbf{U}_{j,j^-}^{(s)}\|_F^2$ , where  $\mathbf{H}_{j,j^-}$  is the  $j$ th row of  $\mathbf{H}$  with  $\mathbf{H}_{jj}$  excluded,  $\mathbf{Z}_{j^-}$  is  $\mathbf{Z}$  with its  $j$ th column removed and  $\mathbf{x}_j$  is the  $j$ th column of  $\mathbf{X}$ . The minimizer is the solution to  $(\mathbf{Z}_{j^-}^{\top} \mathbf{Z}_{j^-} + \rho \mathbf{I}) \mathbf{H}_{j,j^-} = \mathbf{Z}_{j^-}^{\top} \mathbf{z}_j + \rho(\mathbf{F}_{j,j^-}^{(s)} - \mathbf{U}_{j,j^-}^{(s)})$ , where the factorization of  $\mathbf{Z}_{j^-}^{\top} \mathbf{Z}_{j^-} + \rho \mathbf{I}$  can be cached to speed up subsequent updates.

#### A.2.2. C-step

By Lemma 1, the updating formula for  $\mathbf{C}$  is  $\mathbf{C}^{(s+1)} = \boldsymbol{\Lambda}^{1/2} \mathbf{P} \mathbf{O}^{\top}$ , where  $\mathbf{O}$  and  $\mathbf{P}$  are obtained from singular value decomposition of  $(\mathbf{B}^{(s)} - \mathbf{V}^{(s)})^{\top} \boldsymbol{\Lambda}^{1/2} = \mathbf{O} \mathbf{E} \mathbf{P}^{\top}$ .

#### A.2.3. F-step

**F**-step updates two parts of the matrix, one is the adjacency matrix, one is the intervention matrix.

**Part I: adjacency matrix** For  $i = 1, \dots, p$  and  $j = 1, \dots, p$ , we solve the following problem:

$$\min_F \mu_1 \sum_{i,j} |F_{ij}| w_{ij}^{m-1} + \frac{\rho}{2} \sum_{ijk} (|F_{ij}| w_{ij}^{m-1} + \tau(1 - w_{ij}^{m-1}) - L_{ijk}^s) + \frac{\rho}{2} \sum_{i,j} (A_{ij}^{s+1} - F_{ij} + U_{ij}^s)^2,$$

where  $L_{ijk}^s = \lambda_{ik}^s + \tau I(j \neq k) - \lambda_{jk}^s - \xi_{ijk}^s - y_{ijk}^s$ .

We can solve  $F_{ij}$  elementwise:

$$F_{ij}^{s+1} = \begin{cases} S\left(\frac{A_{ij}^{s+1} + U_{ij}^s}{1+p}, \frac{\mu_1 - p\rho \sum_k L_{ijk}^s}{\rho(1+p)}\right) & \text{if } w_{ij}^{m-1} = 1, \\ A_{ij}^{s+1} + U_{ij}^s & \text{if } w_{ij}^{m-1} = 0, \end{cases} \text{ where}$$

$$S(b, \lambda) = \begin{cases} b - 0.5\lambda & \text{if } b > 0.5\lambda, \\ b + 0.5\lambda & \text{if } b < -0.5\lambda, \\ 0 & \text{otherwise,} \end{cases} \text{ is the soft-thresholding operator.}$$

**Part II: intervention matrix** For  $i = 1, \dots, p$  and  $j = p+1, \dots, p+W$ , we solve the following problem:

$$\min_F \mu_2 \sum_{i,j} |F_{ij}| v_{ij}^{m-1} + \frac{\rho}{2} \sum_{i,j} (A_{ij}^{s+1} - F_{ij} + U_{ij}^s)^2 + \frac{\rho}{2} \sum_{i,j} (C_{i,j-p}^s - F_{ij} + Z_{i,j-p}^s)^2.$$

We can solve  $F_{ij}$  elementwise:

$$F_{ij}^{s+1} = \begin{cases} S\left(\frac{1}{2}(A_{ij}^{s+1} + U_{ij}^s + C_{i,j-p}^s + Z_{i,j-p}^s), \frac{\mu_2}{2\rho}\right) & \text{if } v_{ij}^{m-1} = 1, \\ \frac{1}{2}(A_{ij}^{s+1} + U_{ij}^s + C_{i,j-p}^s + Z_{i,j-p}^s) & \text{if } v_{ij}^{m-1} = 0. \end{cases}$$

#### A.2.4. $\lambda$ -step and $\xi$ -step

$(\lambda^{s+1}, \xi^{s+1})$  is updated by  $\lambda^{s+1} = \mathbf{M}_{p \times p} \mathbf{W}_{p \times p}^{s+1}$ , where

$$\mathbf{M}_{p \times p} = \frac{1}{\tau} \begin{pmatrix} 1 & 0 & 0 & \dots & 0 \\ 1 & \frac{2}{p} & \frac{1}{p} & \dots & \frac{1}{p} \\ 1 & \frac{1}{p} & \frac{2}{p} & \dots & \frac{1}{p} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & \frac{1}{p} & \dots & \frac{1}{p} & \frac{2}{p} \end{pmatrix}$$

$$\begin{aligned} W_{1j}^{s+1} &= 1, \\ W_{ik}^{s+1} &= \frac{1}{2}(\tau + \sum_j (|B_{kj}^{s+1}| w_{ij}^{(m-1)} + \tau(1 - w_{ij}^{(m-1)}) + \xi_{ijk}^{s+1} + y_{ijk}^s) \\ &\quad - \sum_j (|B_{kj}^{s+1}| w_{ij}^{(m-1)} + \tau(1 - w_{ij}^{(m-1)}) + \xi_{jik}^{s+1} + y_{jik}^s)); \quad i \neq k, \end{aligned}$$

$$W_{kk}^{s+1} = \frac{1}{2}(-(p-1)\tau + \sum_j (|B_{kj}^{s+1}|w_{ij}^{(m-1)} + \tau(1-w_{ij}^{(m-1)}) + \xi_{kjk}^{s+1} + y_{kjk}^s) - \sum_j (|B_{kj}^{s+1}|w_{ij}^{(m-1)} + \tau(1-w_{ij}^{(m-1)}) + \xi_{jkk}^{s+1} + y_{jkk}^s)),$$

for  $i, j, k = 1, \dots, p$ .

$$\xi_{ijk}^{s+1} := \max(0, (\tau\lambda_{ik}^s + \tau I(j \neq k) - \tau\lambda_{jk}^s - |B_{ij}^{s+1}| - y_{ijk}^s)); i, j, k = 1, \dots, p.$$

### A.3. Computation details for estimating $D_0$

A good estimate of  $D_0$  is obtained by solving (9) without the variance constraint (3). First, we introduce a re-parametrization to exploit the convexity. Let  $\mathbf{R} = \mathbf{D}^{-1/2}$ ,  $\Phi = \mathbf{D}^{-1/2}\mathbf{A}$  and  $\Psi = \mathbf{D}^{-1/2}\mathbf{B}$ , then (2) becomes  $\mathbf{Y} = \mathbf{R}^{-1}\Phi\mathbf{Y} + \mathbf{R}^{-1}\Psi\mathbf{X} + \mathbf{R}^{-1}\epsilon$ ,  $\epsilon \sim N(0, \mathbf{I})$  and  $\Phi = (\phi_{jk})_{p \times p}$ ,  $\Psi = (\psi_{jw})_{p \times W}$  are scaled versions of adjacency matrix  $\mathbf{A}$  and intervention matrix  $\mathbf{B}$  and  $\mathbf{R} = (r_1, \dots, r_p)$ . Under the new parametrization, the likelihood (4) becomes  $l(\Phi, \Psi, \mathbf{R}) = \sum_{j=1}^p \left[ -n \log r_j + \frac{1}{2} \sum_{i=1}^n \left( r_j y_{ij} - \sum_{k \neq j} \phi_{jk} y_{ik} - \sum_{w=1}^W \psi_{jw} x_{iw} \right)^2 \right]$ , an easy check will confirm that  $l(\Phi, \Psi, \mathbf{R})$  is convex in  $(\Phi, \Psi, \mathbf{R})$ . Then (9) without variance constraint (3) is written as

$$\begin{aligned} & \min_{(\Phi, \Psi, \mathbf{R})} l(\Phi, \Psi, \mathbf{R}) = \\ & \sum_{j=1}^p \left[ -n \log r_j + \frac{1}{2} \sum_{i=1}^n \left( r_j y_{ij} - \sum_{k \neq j} \phi_{jk} y_{ik} - \sum_{w=1}^W \psi_{jw} x_{iw} \right)^2 \right], \\ & \text{subj to } \sum_{1 \leq j < l \leq p} J_\tau(\phi_{jl}) \leq K_1, \quad \sum_{1 \leq j \leq p, 1 \leq l \leq W} J_\tau(\psi_{jl}) \leq K_2, \\ & \sum_{j_1=j_{L+1}, 1 \leq k \leq L} J_\tau(\phi_{j_1 k_{j_1+1}}) \leq L-1; \text{ any } (j_1, \dots, j_L), L=2, \dots, p. \end{aligned} \quad (21)$$

To solve (21), following the computation strategy illustrated in Appendix 7.1 and 7.2, after transformation, in the  $m$ th DC step, we solve

$$\begin{aligned} & \min_{(\Phi, \Psi, \mathbf{F}, \lambda, \xi, \mathbf{y}, \mathbf{U}, \mathbf{R})} L_\rho(\Phi, \Psi, \mathbf{F}, \lambda, \xi, \mathbf{y}, \mathbf{U}, \mathbf{R}) = l(\Phi, \Psi, \mathbf{R}) \\ & + \frac{\mu_1}{\tau} \sum_{1 \leq j \neq l \leq p} |F_{jl}| w_{jl}^{(m-1)} + \frac{\mu_2}{\tau} \sum_{1 \leq j \leq p, 1 \leq l \leq W} |\psi_{jl}| v_{jl}^{(m-1)} \\ & + \sum_{1 \leq s \leq p} \sum_{1 \leq j \neq l \leq p} \frac{\rho}{2} \left( |F_{jl}| w_{jl}^{(m-1)} + \tau(1-w_{jl}^{(m-1)}) + \xi_{jls} - \lambda_{jl} \right. \\ & \quad \left. - \tau I(l \neq s) + \lambda_{ls} + y_{jls} \right)^2 \\ & + \frac{\rho}{2} \sum_{1 \leq j, l \leq p} (\phi_{jl} - F_{jl} + U_{jl})^2 \\ & + \frac{\rho}{2} \sum_{1 \leq j \leq p, 1 \leq l \leq W} (\psi_{jl} - F_{j,l+p} + U_{j,l+p})^2, \end{aligned} \quad (22)$$

which is solved over blocks  $(\Phi, \Psi, \mathbf{F}, \lambda, \xi, \mathbf{y}, \mathbf{U}, \mathbf{R})$ . The updating formula are similar to those in Appendix 7.2, with one exception, which is  $\mathbf{R}$ -step, which is illustrated in the following section.

After solving (21), estimates  $(\hat{\Phi}, \hat{\Psi}, \hat{\mathbf{R}})$  are obtained, then  $D_0 = \hat{\mathbf{R}}^{-2}$ .

#### A.3.1. $\mathbf{R}$ -step

Denote  $\mathbf{T} = (\Phi, \Psi)$  be the concatenation of  $\Phi$  and  $\Psi$ , let  $\mathbf{Z} = (\mathbf{Y}, \mathbf{X})$  be the concatenated data matrices we solve the minimization problem:  $\min_{r_j} -n \log r_j +$

$\frac{1}{2} \|r_j \mathbf{z}_j - \mathbf{Z}_j - \mathbf{T}_{jj}^T\|^2$ . The derivative equation becomes  $-\frac{n}{r_j} + \sum_{i=1}^n (r_j z_{ij} - \sum_{k \neq j} T_{jk} z_{ik}) z_{ij} = 0$ , yielding a solution  $r_j = \frac{b + \sqrt{b^2 + 4na}}{2a}$ , where  $a = \sum_{i=1}^n z_{ij}^2$  and  $b = \sum_{i=1}^n z_{ij} \sum_{k \neq j} T_{jk} z_{ik}$ .

#### A.4. Technical proofs

*Proof of Lemma 2.* Let  $\boldsymbol{\epsilon}' = \mathbf{B}\mathbf{X} + \boldsymbol{\epsilon}$ , under (3), (2) reduces to

$$\mathbf{Y} = \mathbf{A}\mathbf{Y} + \boldsymbol{\epsilon}', \quad (23)$$

where  $\boldsymbol{\epsilon}' \sim N(0, \theta \mathbf{I})$ . In (23),  $\boldsymbol{\Omega} = (\mathbf{I} - \mathbf{A})^\top (\mathbf{I} - \mathbf{A}) / \theta$ . Given  $\mathbf{A}$  and the distribution of  $\mathbf{Y}$ , the value of  $\theta$  is unique. This completes the proof.  $\square$

*Proof of Theorem 1.* Let  $\boldsymbol{\epsilon}' = \mathbf{B}\mathbf{X} + \boldsymbol{\epsilon}$ . Note that  $\mathbf{X} \sim N(0, \boldsymbol{\Sigma}_X)$ . Hence  $\boldsymbol{\epsilon}' \sim N(0, \theta \mathbf{I})$ , which is independent of any specific value of  $\mathbf{B}$ . Now (2) is written as

$$\mathbf{Y} = \mathbf{A}\mathbf{Y} + \boldsymbol{\epsilon}'. \quad (24)$$

By (3),  $\boldsymbol{\epsilon}' \sim N(0, \theta \mathbf{I})$ , (24) reduces to an observational model with equal error variance. By Theorem 1 of [18],  $\mathbf{A}$  is identifiable from the distribution of  $\mathbf{Y}$ , which in turn is identifiable from the joint distribution  $(\mathbf{Y}, \mathbf{X})$  that is proportional to the distribution of  $\mathbf{Y}$ . This completes the proof.  $\square$

*Proof of Lemma 1.* The update step for  $\mathbf{C}$  in ADMM is

$$\mathbf{C}^{(k+1)} = \operatorname{argmin}_{\mathbf{C}} \|\mathbf{C} - \mathbf{B}^{(k)} + \mathbf{V}^{(k)}\|_F^2, \text{ s.t. } \mathbf{C}\mathbf{C}^\top = \boldsymbol{\Lambda}, \quad (25)$$

(25) is equivalent to  $\max_{\mathbf{C}} \operatorname{Tr} [\mathbf{C}^\top (\mathbf{B}^{(k)} - \mathbf{V}^{(k)})]$ , s.t.  $\mathbf{C}\mathbf{C}^\top = \boldsymbol{\Lambda}$ , we can transform it into a standard form by introducing  $\mathbf{Q} = \mathbf{C}^\top \boldsymbol{\Lambda}^{-1/2}$  and  $\mathbf{S} = (\mathbf{B}^{(k)} - \mathbf{V}^{(k+1)})^\top \boldsymbol{\Lambda}^{1/2}$ . Then (25) becomes

$$\max_{\mathbf{Q}} \operatorname{Tr} [\mathbf{Q}^\top \mathbf{S}], \text{ s.t. } \mathbf{Q}^\top \mathbf{Q} = \mathbf{I}. \quad (26)$$

To solve (26), we can do a singular value decomposition  $\mathbf{S} = \mathbf{O}\mathbf{E}\mathbf{P}^\top$ , where  $\mathbf{O} \in \mathbb{R}^{w \times p}$ ,  $\mathbf{E}, \mathbf{P} \in \mathbb{R}^{p \times p}$ ,  $\mathbf{O}^\top \mathbf{O} = \mathbf{P}^\top \mathbf{P} = \mathbf{P}\mathbf{P}^\top = \mathbf{I}$ . Then the cost function in (26) can be written in the form  $\operatorname{Tr} [\mathbf{Q}^\top \mathbf{O}\mathbf{E}\mathbf{P}^\top] = \operatorname{Tr} [(\mathbf{Q}\mathbf{P})^\top \mathbf{O}\mathbf{E}] = \sum_{j=1}^p [\mathbf{M}^\top \mathbf{O}]_{jj} E_{jj}$ , where  $\mathbf{M} = \mathbf{Q}\mathbf{P} \in \mathbb{R}^{w \times p}$  and  $[\mathbf{M}^\top \mathbf{O}]_{jj}$  denotes the  $j$ th diagonal element of the cross product  $\mathbf{M}^\top \mathbf{O}$ . Note that the constraint  $\mathbf{Q}^\top \mathbf{Q} = \mathbf{I}$  is equivalent to  $\mathbf{M}^\top \mathbf{M} = \mathbf{I}$  since  $\mathbf{P}$  is an orthogonal matrix. Since  $E_{jj} \geq 0$ , if we can maximize each of  $[\mathbf{M}^\top \mathbf{O}]_{jj}; j = 1, \dots, p$  at the same time, then the problem is solved.

Since  $[\mathbf{M}^\top \mathbf{O}]_{jj} = \sum_{i=1}^W M_{ij} O_{ij}$ , by the Cauchy-Schwartz inequality and the fact that  $\mathbf{M}^\top \mathbf{M} = \mathbf{O}^\top \mathbf{O} = \mathbf{I}$ , we have  $\sum_{i=1}^W M_{ij} O_{ij} \leq \sqrt{\sum_{i=1}^W M_{ij}^2 \sum_{i=1}^W O_{ij}^2} = \sqrt{1 \times 1} = 1$ , where the equality holds if and only if  $M_{ij} = O_{ij}$  for  $i = 1, 2, \dots, W$ . If we maximize each  $[\mathbf{M}^\top \mathbf{O}]_{jj}; j = 1, \dots, p$  at the same time, we need  $\mathbf{M} = \mathbf{O}$ , which leads to  $\mathbf{Q} = \mathbf{O}\mathbf{P}^\top$ . Then the solution of (25) is  $\mathbf{C}^{(k+1)} = \boldsymbol{\Lambda}^{1/2} \mathbf{P}\mathbf{O}^\top$ . This completes the proof.  $\square$

*Proof of Theorem 2.* When we ignore the intervention covariates  $\mathbf{X}$  and pull them into the error terms, under the variance constraint (3), (2) reduces to (24) with equal error variances, under which  $\mathbf{A}$  is identifiable from  $\mathbf{\Omega}$ . Note that  $\mathbf{\Omega}^{obs} = (\mathbf{I} - \mathbf{A}^{obs})^\top (\mathbf{I} - \mathbf{A}^{obs}) / \theta$  is a function of  $(\mathbf{A}^{obs}, \theta)$ , it follows that  $P(\hat{\theta} \neq \hat{\theta}_O) \leq P(\hat{\mathbf{\Omega}}^{obs} \neq \hat{\mathbf{\Omega}}_O^{obs})$ . The rest of the proof follows from Theorem 3 of [27]. This completes the proof.  $\square$

*Proof of Theorem 3.* First, we define a complexity measure for the size of a space  $\mathcal{F}$ . The bracketing Hellinger metric entropy of  $\mathcal{F}$ , denoted by  $H(\cdot, \mathcal{F})$ , is the logarithm of the cardinality of the  $u$ -bracketing of  $\mathcal{F}$  of the smallest size. That is, for a bracket covering  $S(\epsilon, m) = \{f_1^l, f_1^u, \dots, f_m^l, f_m^u\} \subset \mathcal{L}_2$  satisfying  $\max_{1 \leq j \leq m} \|f_j^u - f_j^l\|_2 \leq \epsilon$  and for any  $f \in \mathcal{F}$ , there exists a  $j$  such that  $f_j^l \leq f \leq f_j^u$ , a.e.  $P$ , then  $H(u, \mathcal{F}) = \log(\min\{m : S(u, m)\})$ , where  $\|f\|_2 = \int f^2(z) d\mu$ , with  $\mu$  the dominating measure.

Denote  $E^\tau = \{(i, j) : |A_{ij}| \geq \tau\}$ . When  $K = |E^0|$ ,  $\sum_{1 \leq i \neq j \leq p} J_\tau(A_{ij}) \leq |E^0|$ , so  $|\hat{E}^\tau| \leq |E^0|$ . If  $\hat{E}^\tau = E^0$ , then  $\sum_{1 \leq i \neq j \leq p} |A_{ij}| I(|A_{ij}| < \tau) = 0$ , then  $\hat{\mathbf{A}} = \hat{\mathbf{A}}_O$ . Therefore, it suffices to prove the case when  $\hat{E}^\tau \neq E^0$ .

Define  $\mathbf{\Omega}_{E^\tau} = (\mathbf{I} - \mathbf{A}_{E^\tau})^\top (\mathbf{I} - \mathbf{A}_{E^\tau}) / \theta_{E^\tau}$  for any  $E^\tau \subset \{(i, j) : 1 \leq i \neq j \leq p\}$ . We can partition  $E^\tau$  as  $E^\tau = (E^\tau \setminus E^0) \cap (E \cup E^0)$ . Let  $B_{kj} = \{\mathbf{\Omega}_{E^\tau} : E^\tau \neq E^0, |E^\tau \cap E^0| = k, |E^\tau \setminus E^0| = j, (d_1(|E^0| - k)C_{\min}(\mathbf{\Omega}^0) - d_3 q \tau^{d_2}) \leq h^2(\mathbf{\Omega}_{E^\tau}, \mathbf{\Omega}^0)\}$ ;  $k = 0, \dots, |E^0| - 1, j = 1, \dots, |E^0| - k$ . Then  $B_{kj}$  has  $\binom{|E^0|}{k} \binom{p(p-1)-|E^0|}{j}$  different elements  $E^\tau$ 's of sizes  $|E^\tau \cap E^0| = k, |E^\tau \setminus E^0| = j$ . By definition  $\{\mathbf{\Omega}_{E^\tau} : E^\tau \neq E^0, |E^\tau \cap E^0| \leq |E^0|, C_{\min}(\mathbf{\Omega}^0) \leq h^2(\mathbf{\Omega}_{E^\tau}, \mathbf{\Omega}^0)\} \subset \cup_{k=0}^{|E^0|-1} \cup_{j=1}^{|E^0|-k} B_{kj}$ . Let  $L(\mathbf{\Omega}) = \log f(\mathbf{\Omega}, y, x)$  where  $f(\mathbf{\Omega}, y, x) = f(\mathbf{\Omega}, y|x) f_X(x)$  is the joint density. Then

$$\begin{aligned} & P(\hat{G} \neq G^0) \\ & \leq P(\hat{\mathbf{\Omega}} \neq \hat{\mathbf{\Omega}}_O) \\ & \leq P^* \left( \sup_{\mathbf{\Omega}_{E^\tau} : E^\tau \neq E^0, |E^\tau| \leq |E^0|} (L(\mathbf{\Omega}_{E^\tau}) - L(\hat{\mathbf{\Omega}}_O)) \geq 0 \right) \\ & \leq P^* \left( \sup_{\mathbf{\Omega}_{E^\tau} : E^\tau \neq E^0, |E^\tau| \leq |E^0|} (L(\mathbf{\Omega}_{E^\tau}) - L(\mathbf{\Omega}^0)) \geq 0 \right) \\ & \leq \sum_{E^\tau \subset \{(i,j): 1 \leq i < j \leq p\} : E^\tau \neq E^0, |E^\tau| \leq |E^0|} P^* \left( \sup_{\mathbf{\Omega}_{E^\tau} \in B_{kj}} (L(\mathbf{\Omega}_{E^\tau}) - L(\hat{\mathbf{\Omega}}_O)) \geq 0 \right) \\ & \equiv I, \end{aligned}$$

where  $P^*$  is the outer measure.

For  $I$ , we apply Theorem 1 of [26] to bound each term in the sum. We verify the entropy condition (3.1) there for the bracketing entropy over  $B_{kj}$ . Let  $\mathbf{\Pi}$  denote the covariance matrix of the joint distribution of  $(\mathbf{Y}, \mathbf{X})$  and  $\mathbf{\Gamma} = \mathbf{\Pi}^{-1}$ , then  $\mathbf{\Omega}$  is the upper  $p \times p$  diagonal block of  $\mathbf{\Gamma}$ . Let  $\mathcal{F}_{kj} = \{f^{1/2}(\mathbf{\Gamma}, \cdot, \cdot) : \mathbf{\Omega} \in B_{kj}\}$  be the class of square-root densities. Define  $\mathbf{\Delta} = \tilde{\mathbf{\Gamma}} - \mathbf{\Gamma}$  and let  $\lambda_1, \dots, \lambda_{p+W}$  be the eigenvalues of  $\sqrt{\mathbf{\Pi}} \mathbf{\Delta} \sqrt{\mathbf{\Pi}}$ ,  $z = (y^\top, x^\top)^\top$ . Then  $\max_{1 \leq i \leq p+W} \lambda_i \leq \lambda_{\max}(\mathbf{\Delta}) \times \lambda_{\max}(\mathbf{\Pi}) \leq c(p+W) \|\mathbf{\Pi}\|_{\max}$  following Prob.III.6.14 [1]. By Lemma 6.5 of [28],

it can be shown

$$\begin{aligned} |\log f(\tilde{\Gamma}, y, x) - \log f(\Gamma, y, x)| &\leq \max_{1 \leq i \leq p+W} \lambda_i(D(z) + p + W) \\ &\leq c(p + W) \|\Delta\|_{\max}(D(z) + p + W), \end{aligned}$$

where  $\lambda_{\max}$  denotes the largest eigenvalue,  $D(z) = \lambda_{\max}(\Gamma) \operatorname{tr}(zz^\top) \leq M_4 \times \operatorname{tr}(zz^\top)$ . Note that  $E \operatorname{tr}(zz^\top) \leq c(p + W)$  for some constant  $c > 0$ . By Assumption A.2,  $\lambda_{\max}^2(\sqrt{\Pi}) \leq 1/M_3$ . Then

$$\begin{aligned} I' &\equiv \int \sup_{\tilde{\mathbf{r}} \in B_\delta(\Gamma)} (f^{1/2}(\tilde{\Gamma}, y, x) - f^{1/2}(\Gamma, y, x))^2 d\mu \\ &\leq \sup_{\tilde{\mathbf{r}} \in B_\delta(\Gamma)} c'(p + W) \|\Delta\|_{\max}^2 E(D(z) + p + W)^2, \end{aligned}$$

for some constant  $c' > 0$ .

Then for some positive constant  $Q$ ,  $I' \leq Q(p + W)^4 \delta^2$ . By Lemma 1 of [15], it suffices to bound the entropy of  $B_{ij}$ . There are  $|E|$  nonzero entries of  $\mathbf{A}$  with  $\binom{p(p-1)}{|E|}$  possible locations. By [13], for  $u \geq \epsilon^2$ ,

$$\begin{aligned} H(u, \mathcal{F}_{ij}) &\leq c_0 \left( \log \binom{p(p-1)}{|E|} + |E| \log \left( \frac{\min(M_2^{1/2}, 1)}{u} \right) \right) \\ &\leq c_0 \left( |E| \log \left( e \frac{p(p-1)}{|E|} \right) + |E| \log \left( \frac{\min(M_2^{1/2}, 1)}{u} \right) \right) \\ &\leq c_0 (|E| \log p \log(1/u)) \end{aligned}$$

Then  $\epsilon = \epsilon_{n,p,|E^0|} = \min(1, (2c_0)^{1/2} c_4^{-1} \log(2^{1/2}/c_3) \log p (|E^0|/n)^{1/2})$  satisfies

$$\sup_{0 \leq |E^\tau| \leq |E^0|} \int_{2^{-8}\epsilon^2}^{2^{1/2}\epsilon} H(t/c_3, \mathcal{F}_{ij}) dt \leq (2|E^0|)^{1/2} \epsilon \log(2^{1/2}/c_3) \leq c_4 n^{1/2} \epsilon^2. \quad (27)$$

for some constants  $c_3, c_4 > 0$ , say  $c_3 = 10, c_4 = \frac{(2/3)^{5/2}}{512}$ . By Assumption B,  $C_{\min}(\Omega^0) \geq \epsilon_{n,p,|E^0|}$  implies (27), provided that  $2d_0^{-1} > (2c_0)^{1/2} c_4^{-1} \log(2^{1/2}/c_3)$ . Using the facts about binomial coefficients:  $\sum_{j=1}^{|E^0|-k} \binom{p(p-1)-|E^0|}{j} \leq (p(p-1) - |E^0| + 1)^{|E^0|-k}$  and  $\binom{|E^0|}{k} \leq |E^0|^k$ , we have, by Theorem 1 of [26], that for a constant  $c_2 > 0$ , say  $c_2 = \frac{4}{27} \frac{1}{1926}$ ,

$$\begin{aligned} I &\leq \sum_{k=0}^{|E^0|-1} \sum_{j=1}^{|E^0|-k} P^* \left( \sup_{\Omega_{E^\tau} \in B_{k,j}} (L(\Omega_{E^\tau}) - L(\hat{\Omega}_O)) \geq 0 \right) \\ &\leq 4 \sum_{k=0}^{|E^0|-1} \binom{|E^0|}{k} \exp(-c_2 n (d_1 C_{\min}(\Omega^0) - d_3 q \tau^{d_2})) \sum_{j=1}^{|E^0|-k} \binom{p(p-1) - |E^0|}{j} \end{aligned}$$

$$\begin{aligned} &\leq 4 \sum_{i=1}^{|E^0|} \exp(-i((c_2 d_1/2)C_{\min}(\mathbf{\Omega}^0) - \log(p(p-1) - |E^0| + 1) - \log |E^0|)) \\ &\leq R(\exp(-(c_2 d_1/2)C_{\min}(\mathbf{\Omega}^0) - \log(p(p-1) - |E^0| + 1) - \log |E^0|)), \end{aligned}$$

provided that  $\tau \leq C_{\min}(\mathbf{\Omega}^0)M_1/4p$ , where  $R(x) = x/(1-x)$ . Moreover, since  $I \leq 1$  and  $\log(p(p-1) - |E^0| + 1) + \log |E^0| \leq 2 \log((p(p-1) + 1)/2)$ ,

$$\begin{aligned} I &\leq 5 \exp(-(c_2 d_1/2)nC_{\min}(\mathbf{\Omega}^0) + 2 \log((p(p-1) + 1)/2)) \\ &\leq \exp(-c_2 n C_{\min} + 2 \log(p(p-1) + 1) + 3). \end{aligned}$$

Under Assumption B,  $P(\hat{G} \neq G^0) \rightarrow 0$  as  $n, p, |E^0| \rightarrow \infty$ . For parameter estimation,

$$\begin{aligned} Eh^2(\hat{\mathbf{\Omega}}, \mathbf{\Omega}^0) &\leq E[h^2(\hat{\mathbf{\Omega}}_O, \mathbf{\Omega}^0)I(\hat{\mathbf{\Omega}} = \hat{\mathbf{\Omega}}_O)] + P(\hat{\mathbf{\Omega}} \neq \hat{\mathbf{\Omega}}_O) \\ &\leq (1 + o(1))Eh^2(\hat{\mathbf{\Omega}}_O, \mathbf{\Omega}^0), \end{aligned}$$

then  $\frac{Eh^2(\hat{\mathbf{\Omega}}, \mathbf{\Omega}^0)}{Eh^2(\hat{\mathbf{\Omega}}_O, \mathbf{\Omega}^0)} \rightarrow 1$  as  $n, p, |E^0| \rightarrow \infty$ .

This completes the proof.  $\square$

## References

- [1] BHATIA, R. (2013). *Matrix analysis* **169**. Springer Science & Business Media. [MR1477662](#)
- [2] BOYD, S., PARIKH, N., CHU, E., PELEATO, B. and ECKSTEIN, J. (2011). Distributed optimization and statistical learning via the alternating direction method of multipliers. *Foundations and Trends in Machine Learning* **3** 1–122.
- [3] CHICKERING, D. M. (2003). Optimal structure identification with greedy search. *The Journal of Machine Learning Research* **3** 507–554. [MR1991085](#)
- [4] EATON, D. and MURPHY, K. P. (2007). Exact Bayesian structure learning from uncertain interventions. In *International Conference on Artificial Intelligence and Statistics* 107–114.
- [5] EDWARDS, D. (2000). *Introduction to graphical modelling*. Springer Verlag. [MR1880319](#)
- [6] ELLIS, B. and WONG, W. H. (2008). Learning causal Bayesian network structures from experimental data. *Journal of the American Statistical Association* **103** 778–789. [MR2524009](#)
- [7] FRIEDMAN, N. (2004). Inferring cellular networks using probabilistic graphical models. *Science Signalling* **303** 799.
- [8] FU, F. and ZHOU, Q. (2013). Learning sparse causal Gaussian networks with experimental intervention: regularization and coordinate descent. *Journal of the American Statistical Association* **108** 288–300. [MR3174620](#)
- [9] HAUSER, A. and BÜHLMANN, P. (2012). Characterization and greedy learning of interventional Markov equivalence classes of directed acyclic graphs. *The Journal of Machine Learning Research* **13** 2409–2464. [MR2973606](#)

- [10] HUANG, S., LI, J., YE, J., FLEISHER, A., CHEN, K., WU, T. and REIMAN, E. (2013). A Sparse Structure Learning Algorithm for Gaussian Bayesian Network Identification from High-Dimensional Data. *Pattern Analysis and Machine Intelligence, IEEE Transactions on* **35** 1328–1342.
- [11] KALISCH, M. and BÜHLMANN, P. (2007). Estimating high-dimensional directed acyclic graphs with the PC-algorithm. *The Journal of Machine Learning Research* **8** 613–636.
- [12] KANEHISA, M. and GOTO, S. (2000). KEGG: kyoto encyclopedia of genes and genomes. *Nucleic acids research* **28** 27–30.
- [13] KOLMOGOROV, A. N. and TIKHOMIROV, V. M. (1959).  $\varepsilon$ -entropy and  $\varepsilon$ -capacity of sets in function spaces. *Uspekhi Matematicheskikh Nauk* **14** 3–86. [MR0112032](#)
- [14] NANDY, P., HAUSER, A., MAATHUIS, M. H. et al. (2018). High-dimensional consistency in score-based and hybrid structure learning. *The Annals of Statistics* **46** 3151–3183. [MR3851768](#)
- [15] OSSIANDER, M. (1987). A central limit theorem under metric entropy with  $L_2$  bracketing. *The Annals of Probability* **15** 897–919. [MR0893905](#)
- [16] PEARL, J. (2003). Statistics and causal inference: A review. *Test* **12** 281–345. [MR2044313](#)
- [17] PENG, S., SHEN, X. and PAN, W. (2019). intdag: Reconstruction of a Directed Acyclic Graph with Interventions R package version 1.0.1.
- [18] PETERS, J. and BÜHLMANN, P. (2013). Identifiability of Gaussian structural equation models with equal error variances. *Biometrika, first published online*, doi: [10.1093/biomet/ast043](https://doi.org/10.1093/biomet/ast043). [MR3180667](#)
- [19] SACHS, K., PEREZ, O., PE'ER, D., LAUFFENBURGER, D. A. and NOLAN, G. P. (2005). Causal protein-signaling networks derived from multiparameter single-cell data. *Science* **308** 523.
- [20] SHEN, X., PAN, W. and ZHU, Y. (2012). Likelihood-based selection and sharp parameter estimation. *Journal of the American Statistical Association* **107** 223–232. [MR2949354](#)
- [21] SHEN, X., PAN, W., ZHU, Y. and ZHOU, H. (2013). On constrained and regularized high-dimensional regression. *Annals of the Institute of Statistical Mathematics* 1–26. [MR3105798](#)
- [22] SPIRITES, P., GLYMOUR, C. N. and SCHEINES, R. (2000). *Causation, prediction, and search* **81**. The MIT Press. [MR1815675](#)
- [23] TRIPATHI, S., CHRISTIE, K. R., BALAKRISHNAN, R., HUNTLEY, R., HILL, D. P., THOMMESEN, L., BLAKE, J. A., KUIPER, M. and LÆGREID, A. (2013). Gene Ontology annotation of sequence-specific DNA binding transcription factors: setting the stage for a large-scale curation effort. *Database* **2013**.
- [24] TSAMARDINOS, I., BROWN, L. E. and ALIFERIS, C. F. (2006). The max-min hill-climbing Bayesian network structure learning algorithm. *Machine Learning* **65** 31–78.
- [25] WEBSTER, J. A., GIBBS, J. R., CLARKE, J., RAY, M., ZHANG, W., HOLMANS, P., ROHRER, K., ZHAO, A., MARLOWE, L., KALEEM, M. et al. (2009). Genetic control of human brain transcript expression in Alzheimer

- disease. *The American Journal of Human Genetics* **84** 445–458.
- [26] WONG, W. H. and SHEN, X. (1995). Probability inequalities for likelihood ratios and convergence rates of sieve MLEs. *The Annals of Statistics* **23** 339–362. [MR1332570](#)
- [27] YUAN, Y., SHEN, X., PAN, W. and WANG, Z. (2018). Constrained likelihood for reconstructing a directed acyclic Gaussian graph. *Biometrika* **106** 109–125. [MR3912386](#)
- [28] ZHOU, S., SHEN, X., WOLFE, D. A. et al. (1998). Local asymptotics for regression splines and confidence regions. *The annals of statistics* **26** 1760–1782. [MR1673277](#)