# ParceLiNGAM: A causal ordering method robust against latent confounders

Tatsuya Tashiro[*], Shohei Shimizu[†], Aapo Hyvärinen[‡] and Takashi Washio[§]

## Abstract

We consider learning a causal ordering of variables in a linear non-Gaussian acyclic model called LiNGAM. Several existing methods have been shown to consistently estimate a causal ordering assuming that all the model assumptions are correct. But, the estimation results could be distorted if some assumptions actually are violated. In this paper, we propose a new algorithm for learning causal orders that is robust against one typical violation of the model assumptions: latent confounders. The key idea is to detect latent confounders by testing independence between estimated external influences and find subsets (parcels) that include variables that are not affected by latent confounders. We demonstrate the effectiveness of our method using artificial data and simulated brain imaging data.

## 1 Introduction

Bayesian networks have been widely used to analyze causal relations of variables in many empirical sciences (Bollen, 1989; Pearl, 2000; Spirtes, Glymour, & Scheines, 1993). A common assumption is linear-Gaussianity. But this poses serious identifiability problems so that many important models are indistinguishable with no prior knowledge on the structures. Recently, it was shown by (Shimizu, Hoyer, Hyvärinen, & Kerminen, 2006) that the utilization of non-Gaussianity allows the full structure of a linear acyclic model to be identified without pre-specifying any causal orders of variables. The new model, a Linear Non-Gaussian Acyclic Model called LiNGAM (Shimizu et al., 2006), is closely related to independent component analysis (ICA) (Hyvärinen, Karhunen, & Oja, 2001).

Most existing estimation methods (Shimizu et al., 2006, 2011; Hyvärinen & Smith, 2013) for LiNGAM learn causal orders assuming that all the model assumptions

---

[*]The Institute of Scientific and Industrial Research (ISIR), Osaka University, Mihogaoka 8-1, Ibaraki, Osaka 567-0047, Japan. Email: tashiro@ar.sanken.osaka-u.ac.jp

[†]Osaka University, Japan

[‡]University of Helsinki, Finland

[§]Osaka University, Japan

hold. Therefore, these algorithms could return completely wrong estimation results when some of the model assumptions are violated. Thus, in this paper, we propose a new algorithm for learning causal orders that is robust against one typical model violation, *i.e.*, latent confounders. A latent confounder means a variable which is not observed but which exerts a causal influence on some of the observed variables. Many real-world applications including brain imaging data analysis (Smith et al., 2011) could benefit from our approach.

This paper[1] is organized as follows. We first review LiNGAM (Shimizu et al., 2006) and its extension to latent confounder cases (Hoyer, Shimizu, Kerminen, & Palviainen, 2008) in Section 2. In Section 3, we propose a new algorithm to learn causal orders in LiNGAM with latent confounders. We empirically evaluate the performance of our algorithm using artificial data in Section 4 and simulated fMRI data in Section 5. We conclude this paper in Section 6.

## 2 Background: LiNGAM with latent confounders

We briefly review a linear non-Gaussian acyclic model called LiNGAM (Shimizu et al., 2006) and an extension of the LiNGAM to cases with latent confounding variables (Hoyer et al., 2008).

In LiNGAM (Shimizu et al., 2006), causal relations of observed variables $x_i$ $(i = 1, \cdots, d)$ are modeled as:

$$x_i \quad = \quad \sum_{k(j) < k(i)} b_{ij} x_j + e_i, \tag{1}$$

where $k(i)$ is a causal ordering of the variables $x_i$. In this ordering, the variables $x_i$ graphically form a directed acyclic graph (DAG) so that no later variable determines, *i.e.*, has a directed path on any earlier variable. The $e_i$ are external influences, and $b_{ij}$ are connection strengths. In matrix form, the model (1) is written as

$$\boldsymbol{x} \quad = \quad \mathbf{B}\boldsymbol{x} + \boldsymbol{e}, \tag{2}$$

where the connection strength matrix $\mathbf{B}$ collects $b_{ij}$ and the vectors $\boldsymbol{x}$ and $\boldsymbol{e}$ collect $x_i$ and $e_i$. Note that the matrix $\mathbf{B}$ can be permuted to be lower triangular with all zeros on the diagonal if simultaneous equal row and column permutations are made according to a causal ordering $k(i)$ because of the acyclicity. The zero/non-zero pattern of $b_{ij}$ corresponds to the absence/existence pattern of directed edges. External influences $e_i$ follow non-Gaussian continuous distributions with zero mean and non-zero variance and are mutually independent. The non-Gaussianity assumption on $e_i$ enables identification of a causal ordering $k(i)$ based on data $\boldsymbol{x}$ only (Shimizu et al., 2006). This feature is a major advantage over conventional Bayesian networks based on the Gaussianity assumption on $e_i$ (Spirtes et al., 1993).

---

[1]Some preliminary results were presented in (Tashiro, Shimizu, Hyvärinen, & Washio, 2012), which corresponds to Section 3.1 of this paper.

Next, LiNGAM with latent confounders (Hoyer et al., 2008) can be formulated as follows:

$$\boldsymbol{x} \;\; = \;\; \mathbf{B}\boldsymbol{x} + \boldsymbol{\Lambda}\boldsymbol{f} + \boldsymbol{e}, \tag{3}$$

where the difference with LiNGAM in Eq. (2) is the existence of latent confounding variable vector $\boldsymbol{f}$. A latent confounding variable is a latent variable that is a parent of more than one observed variable. The vector $\boldsymbol{f}$ collects non-Gaussian latent confounders $f_j$ with zero mean and non-zero variance $(j = 1, \cdots, q)$. Without loss of generality (Hoyer et al., 2008), latent confounders $f_j$ are assumed to be mutually independent. The matrix $\boldsymbol{\Lambda}$ collects $\lambda_{ij}$ which denotes the connection strength from $f_j$ to $x_i$. For each $j$, at least two $\lambda_{ij}$ are non-zero since a latent confounder is defined to have at least two children (Hoyer et al., 2008). The matrix $\boldsymbol{\Lambda}$ is assumed to be of full column rank.

The central problem of causal discovery based on the latent variable LiNGAM in Eq. (3) is to estimate *as many* of causal orders $k(i)$ and connection strengths $b_{ij}$ *as possible* based on data $\boldsymbol{x}$ only. This is because in many cases only an equivalence class of the true model whose members produce the exact same observed distribution is identifiable (Hoyer et al., 2008).

In (Hoyer et al., 2008), an estimation method based on overcomplete ICA (Lewicki. & Sejnowski, 2000) was proposed. However, overcomplete ICA methods are often not very reliable and get stuck in local optima. Thus, in (Entner & Hoyer, 2011), a method that does not use overcomplete ICA was proposed to first find variable *pairs* that are not affected by latent confounders and then estimate a causal ordering of one to the other. However, their method does not estimate a causal ordering of more than two variables. A simple cumulant-based method for estimating the model in the case of Gaussian latent confounders was further proposed by (Chen & Chan, 2013).

# 3   A method robust against latent confounders

In this section, we propose a new approach for estimating causal orders of more than two variables without explicitly modeling latent confounders.

## 3.1   Identification of causal orders of variables that are *not* affected by latent confounders

We first provide principles to identify an exogenous (root) variable and a sink variable which are such that are not affected by latent confounders in the latent variable LiNGAM in Eq. (3) (if such variables exist) and next present an estimation algorithm. Recent estimation methods (Shimizu et al., 2011) for LiNGAM in Eq. (2) and its nonlinear extension (Hoyer, Janzing, Mooij, Peters, & Schölkopf, 2009; Mooij, Janzing, Peters, & Schölkopf, 2009) learn a causal ordering by finding causal orders one by one either from the top downward or from the bottom upward assuming no latent confounders. We extend these ideas to latent confounder cases.

We first generalize Lemma 1 of (Shimizu et al., 2011) for the case of latent confounders.

**Lemma 1** *Assume that all the model assumptions of the latent variable LiNGAM in Eq. (3) are met and the sample size is infinite. Denote by $r_i^{(j)}$ the residuals when $x_i$ are regressed on $x_j$: $r_i^{(j)} = x_i - \frac{\text{cov}(x_i, x_j)}{\text{var}(x_j)} x_j$ $(i \neq j)$. Then a variable $x_j$ is an exogenous variable in the sense that it has no parent observed variable nor latent confounder if and only if $x_j$ is independent of its residuals $r_i^{(j)}$ for all $i \neq j$.* $\square$

Next, we generalize the idea of (Mooij et al., 2009) for the case of latent confounders.

**Lemma 2** *Assume that all the model assumptions of the latent variable LiNGAM in Eq. (3) are met and the sample size is infinite. Denote by $\boldsymbol{x}_{(-j)}$ a vector that contains all the variables other than $x_j$. Denote by $r_j^{(-j)}$ the residual when $x_j$ is regressed on $\boldsymbol{x}_{(-j)}$, i.e., $r_j^{(-j)} = x_j - \boldsymbol{\sigma}_{(-j)j}^T \Sigma_{(-j)}^{-1} \boldsymbol{x}_{(-j)}$, where $\Sigma = \begin{bmatrix} \sigma_j & \boldsymbol{\sigma}_{j(-j)}^T \\ \boldsymbol{\sigma}_{j(-j)} & \Sigma_{(-j)} \end{bmatrix}$ is the covariance matrix of $[x_j, \boldsymbol{x}_{(-j)}^T]^T$. Then a variable $x_j$ is a sink variable in the sense that it has no child observed variable nor latent confounder if and only if $\boldsymbol{x}_{(-j)}$ is independent of its residual $r_j^{(-j)}$.* $\square$

The proofs of these lemmas are given in the appendix.[2]

Thus, we can take a hybrid estimation approach that uses these two principles. We first identify an exogenous variable by finding a variable that is most independent of its residuals and remove the effect of the exogenous variable from the other variables by regressing it out. We repeat this until independence between every variable and all of its residuals is statistically rejected. Dependency between every variable and any of its residuals implies that an exogenous variable as defined in Lemma 1 does not exist or some model assumption of latent variable LiNGAM in Eq. (3) is violated. Similarly, we next identify a sink variable in the remaining variables by finding a variable such that its regressors and its residual are most independent and disregard the sink variable. We repeat this until independence is statistically rejected for every variable.[3] To test independence, we first evaluate pairwise independence between variables and the residuals using a kernel-based independence measure called HSIC (Gretton et al., 2008) and then combine the resulting $p$-values $p_i$ $(i = 1, \cdots, c)$ using a well-known Fisher's method (Fisher, 1950) to compute the test statistic $-2 \sum_{i=1}^{c} \log p_i$, which follows the chi-square distribution with $2c$ degrees of freedom when all the pairs are independent.

Since all the causal orders are not necessarily identifiable in the latent variable LiNGAM in Eq. (3) (Hoyer et al., 2008), we here aim to estimate a $d \times d$

---

[2]We prove the lemmas without assuming the faithfulness (Spirtes et al., 1993) unlike our previous work (Tashiro et al., 2012).

[3]The issue of multiple comparisons arises in this context, which we would like to study in future work.

causal ordering matrix $\mathbf{C}=[c_{ij}]$ that collects causal orderings between two variables, which is defined as

$$
c_{ij} \quad := \quad \begin{cases} -1 & \text{if } k(i) < k(j) \\ 1 & \text{if } k(i) > k(j) \\ 0 & \text{if it is unknown whether either of the two cases} \\ & \text{above } (-1 \text{ or } 1) \text{ is true.} \end{cases} \tag{4}
$$

Thus, the estimation consists of the following steps:

---

**Algorithm 1: Hybrid estimation of causal orders of variables that are not affected by latent confounders**

---

**INPUT:** Data matrix $\mathbf{X}$ and a threshold $\alpha$

1. Given a $d$-dimensional random vector $\boldsymbol{x}$, a $d \times n$ data matrix of the random vector as $\mathbf{X}$ and a significance level $\alpha$, define $U$ as the set of variable indices of $\boldsymbol{x}$, i.e., $\{1, \cdots, d\}$ and initialize an ordered list of variables $K_{head} := \emptyset$ and $K_{tail} := \emptyset$ and $m := 1$. $K_{head}$ and $K_{tail}$ denote the first $|K_{head}|$ variable indices and the last $|K_{tail}|$ variable indices respectively, where each of $|K_{head}|$ and $|K_{tail}|$ denotes the number of elements in the list.

2. Let $\tilde{\boldsymbol{x}} := \boldsymbol{x}$ and $\tilde{\mathbf{X}} := \mathbf{X}$ and find causal orders one by one from the top downward:

   (a) Do the following steps for all $j \in U \setminus K_{head}$: Perform least squares regressions of $\tilde{x}_i$ on $\tilde{x}_j$ for all $i \in U \setminus K_{head}$ ($i \neq j$) and compute the residual vectors $\tilde{\boldsymbol{r}}^{(j)}$ and the residual matrix $\tilde{\mathbf{R}}^{(j)}$. Then, find a variable $\tilde{x}_m$ that is most independent of its residuals:

   $$
   \tilde{x}_m = \arg \max_{j \in U \setminus K_{head}} P_{Fisher}(\tilde{x}_j, \tilde{\boldsymbol{r}}^{(j)}), \tag{5}
   $$

   where $P_{Fisher}(\tilde{x}_j, \tilde{\boldsymbol{r}}^{(j)})$ is the $p$-value of the test statistic defined as $-2 \sum_i \log\{P_H(\tilde{x}_j, \tilde{r}_i^{(j)})\}$, where $P_H(\tilde{x}_j, \tilde{r}_i^{(j)})$ is the $p$-value of the HSIC.

   (b) Go to Step 3 if $P_{Fisher}(\tilde{x}_m, \tilde{\boldsymbol{r}}^{(m)}) < \alpha$, i.e., all independencies are rejected.

   (c) Append $m$ to the end of $K_{head}$ and let $\tilde{\boldsymbol{x}} := \tilde{\boldsymbol{r}}^{(m)}$ and $\tilde{\mathbf{X}} := \tilde{\mathbf{R}}^{(m)}$. If $|K_{head}| = d - 1$, append the remaining variable index to the end of $K_{head}$ and terminate. Otherwise, go back to Step (2a).

3. If $|K_{head}| < d - 2$, let $\boldsymbol{x}' = \boldsymbol{x}$ and $\mathbf{X}' = \mathbf{X}$ and $U' := U \setminus K_{head}$ and find causal orders one by one from the bottom upward [4]:

---

[4]We do not examine remaining two variables in this step since it is already implied in Step 2 that some latent confounders exist. If there were no latent confounders between the remaining two, their causal orders would have already been estimated in Step 2.

(a) Do the following steps for all $j \in U' \setminus K_{tail}$: Collect all the variables except $x'_j$ in a vector $\boldsymbol{x}'_{(-j)}$. Perform least squares regressions of $x'_j$ on $\boldsymbol{x}'_{(-j)}$ and compute the residual $r'^{(-j)}_j$. Then, find such a variable $x'_m$ that its regressors and its residual are most independent:

$$x'_m = \arg \max_{j \in U' \setminus K_{tail}} P_{Fisher}(\boldsymbol{x}'_{(-j)}, r'^{(-j)}_j). \tag{6}$$

(b) Terminate if $P_{Fisher}(\boldsymbol{x}'_{(-m)}, r'^{(-m)}_m) < \alpha$, $i.e.$, all independencies are rejected.

(c) Append $m$ to the top of $K_{tail}$ and let $\boldsymbol{x}' = \boldsymbol{x}'_{(-m)} \mathbf{X}' = \mathbf{X}'_{(-m)}$. Terminate [4] if $|U' \setminus K_{tail}| < 3$ and otherwise go back to Step (3a).

4. Estimate a causal ordering matrix $\mathbf{C}$ based on $K_{head}$ and $K_{tail}$ as follows. Estimate $c_{ij}$ by -1, $i.e.$, $k(i) < k(j)$ in either of the following cases: i) $i$ is earlier than $j$ in $K_{head}$; ii) $i$ is earlier than $j$ in $K_{tail}$; iii) $i$ is in $K_{head}$ and $j$ is in $K_{tail}$; iv) $i$ is in $K_{head}$ and $j$ is neither $K_{head}$ nor $K_{tail}$. Estimate $c_{ij}$ by 1, $i.e.$, $k(i) > k(j)$ in either of the following cases: i) $i$ is later than $j$ in $K_{head}$; ii) $i$ is later than $j$ in $K_{tail}$; iii) $i$ is in $K_{tail}$ and $j$ is in $K_{head}$; iv) $i$ is in $K_{tail}$ and $j$ is neither $K_{head}$ nor $K_{tail}$. Estimate $c_{ij}$ by 0, $i.e.$, the ordering is unknown if $i$ and $j$ are neither in $K_{tail}$ nor $K_{head}$. Note that causal orders of variables that are not in $K_{head}$ or $K_{tail}$ are no later than any in $K_{tail}$ and no earlier than any in $K_{head}$.

**OUTPUT:** Ordered lists $K_{head}$ and $K_{tail}$ and a causal ordering matrix $\mathbf{C}$

## 3.2 A new estimation algorithm robust against latent confounders

Algorithm 1 outputs no causal orders in cases where exogenous variables and sink variables as in Lemmas 1 and 2 do not exist. For example, in the left of Fig. 1, there is no such exogenous variable or sink variable that is not affected by any latent confounder since the latent confounder $f_1$ affects the exogenous variable $x_1$ and the sink variable $x_4$. Therefore, Algorithm 1 would not find any causal orders. However, if we omit $x_4$ as in the right of Fig. 1 and apply Algorithm 1 on the remaining $x_1, x_2, x_3$ only, it will find all the causal orders of $x_1, x_2, x_3$ since $f_1$ does not affect any two of $x_1, x_2, x_3$ and is no longer a latent confounder. The same idea applies to the case that $x_1$ is omitted.

Thus, we propose applying Algorithm 1 on every subset of variables with the size larger than one. This enables learning more causal orders than analyzing the whole set of variables if a subset of variables has exogenous variables or sink variables that are not affected by latent confounders. In practice, Algorithm 1 could give inconsistent causal orderings between a pair of variables for different subsets of variables because of estimation errors. To manage possible inconsistencies in the many causal orderings thus estimated, we rank the obtained
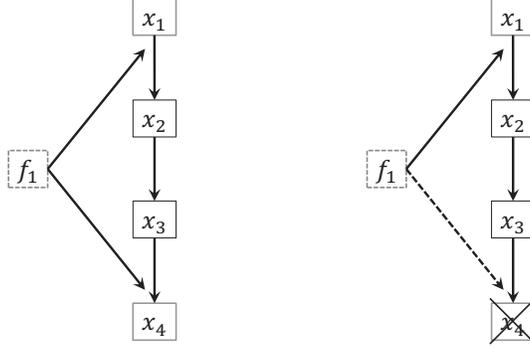
Figure 1: Left: An example graph where Algorithm 1 finds no causal orders. The $f_1$ is a latent confounder that affects $x_1$ and $x_4$. Right: Algorithm 1 finds the causal orders of $x_1$, $x_2$ and $x_3$ if $x_4$ is omitted and only $x_1$, $x_2$ and $x_3$ are analyzed.

causal ordering matrices by plausibility based on the statistical significances (this will be defined below). Then, considering any pair of two variables, we use the causal ordering given by the causal ordering matrix which has the highest plausibility and does contain an estimated causal ordering (*i.e.*, the ordering was not considered unknown) between those two variables.

We evaluate the plausibility of every causal ordering matrix by the $p$-value of the test statistic created based on Fisher's method combining all the $p$-values computed to estimate the causal orders $K_{head}$ and $K_{tail}$ in Algorithm 1. A higher $p$-value can be considered to be more plausible. The test statistic is computed based on $\mathbf{X}$, $K_{head}$ and $K_{tail}$ as follows:

$$-2(\sum_{m \in K_{head}} \sum_{i:k(i)>k(m)} \log\{P_H(\tilde{x}_m, \tilde{r}_i^{(m)})\} + \sum_{m \in K_{tail}} \sum_{i:k(i)<k(m)} \log\{P_H(x_i', r'^{(-m)}_m)\}), \quad (7)$$

where $P_H(\tilde{x}_m, \tilde{r}_i^{(m)})$ and $P_H(x_i', r'^{(-m)}_m)$ are the $p$-values computed to estimate ordered lists $K_{head}$ and $K_{tail}$ in Algorithm 1.

Thus, the estimation consists of the following steps:

---

Algorithm 2: Applying Algorithm 1 on every subset of variables and merging results

---

**INPUT:** Data matrix $\mathbf{X}$ and a threshold $\alpha$

1. Take all the $l$-combinations of variable indices $\{1, \cdots, d\}$ for $l = 2, \cdots, d$. Denote the subsets of variable indices by $U_{subset}^{(s)}$ ($s = 1, \cdots, S$) and the

7

corresponding data matrices by $\mathbf{X}^{(s)}_{subset}$ ($s = 1, \cdots, S$), where $S$ is the number of the subsets.

2. Apply Algorithm 1 on $\mathbf{X}^{(s)}_{subset}$ using the threshold $\alpha$ to estimate $K^{(s)}_{head}$, $K^{(s)}_{tail}$ and $\mathbf{C}^{(s)}$ for all $s \in \{1, \cdots, S\}$, where $K^{(s)}_{head}$ and $K^{(s)}_{tail}$ are ordered lists of Subset $U^{(s)}_{subset}$ and $\mathbf{C}^{(s)}$ is a causal ordering matrix of Subset $U^{(s)}_{subset}$.

3. Compute the $p$-value of the test statistic in Eq. (7) to evaluate the plausibility of $\mathbf{C}^{(s)}$ for all $s \in \{1, \cdots, S\}$.

4. Estimate every element $c_{ij}$ ($i \neq j$) of a causal ordering matrix $\mathbf{C}$ by the causal ordering between $x_i$ and $x_j$ of the causal ordering matrix that has the highest plausibility and does contain an estimated causal ordering between $x_i$ and $x_j$, that is, $k(i) < k(j)$ or $k(j) < k(i)$.

**OUTPUT:** A causal ordering matrix $\mathbf{C}$

---

Algorithm 2 is a brute force approach since it applies Algorithm 1 on *every* subset (parcel) of variables. We could alleviate the computational load by first applying Algorithm 1 on the whole set of variables and then applying Algorithm 2 on the remaining variables whose causal orders have not been estimated after the effects of estimated exogenous variables are removed by regression. Thus, we finally propose the following algorithm called ParceLiNGAM:

---

Algorithm 3: The ParceLiNGAM algorithm

---

**INPUT:** Data matrix $\mathbf{X}$ and a threshold $\alpha$

1. Given a $d$-dimensional random vector $\boldsymbol{x}$ and a $d \times n$ data matrix of the random vector as $\mathbf{X}$, define $U$ as the set of variable indices of $\boldsymbol{x}$, *i.e.*, $\{1, \cdots, d\}$. initialize a $d \times d$ causal ordering matrix $\mathbf{C}$ by the zero matrix.

2. Apply Algorithm 1 on $\mathbf{X}$ using the threshold $\alpha$ to estimate $K_{head}$ and $K_{tail}$ and update $\mathbf{C}$.

3. Let $U_{res} := U \setminus (K_{head} \bigcup K_{tail})$. Denote by $\mathbf{C}_{res}$ the corresponding causal ordering matrix. Denote by $|U_{res}|$ the number of elements in $U_{res}$. Go to Step 6 if $|U_{res}| \leq 2$.

4. Collect variables $x_j$ with $j \in U_{res}$ in a vector $\boldsymbol{x}_{res}$. Collect variables $x_j$ with $j \in K_{head}$ in a vector $\boldsymbol{x}_{head}$. Perform least squares regressions of $\boldsymbol{x}_{head}$ on the $i$-th element of $\boldsymbol{x}_{res}$ for all $i \in U_{res}$ and collect the residuals in the residual matrix $\mathbf{R}_{res}$ whose $i$-th row is given by the residuals regressed on $x_i$.

5. Apply Algorithm 2 on $\mathbf{R}_{res}$ using the threshold $\alpha$ to estimate $\mathbf{C}_{res}$. Replace every $c_{ij}$ ($i \neq j$) of $\mathbf{C}$ by the corresponding element of $\mathbf{C}_{res}$ if $c_{ij}$ is zero and the corresponding element of $\mathbf{C}_{res}$ is 1 or -1.

8

6. Estimate connection strengths $b_{ij}$ if all the non-descendants of $x_i$ are estimated, $i.e.$, the $i$-th row of $\mathbf{C}$ has no zero. This can be done by doing multiple regression of $x_i$ on all of its non-descendants $x_j$ with $k(j) < k(i)$.

**OUTPUT:** A causal ordering matrix $\mathbf{C}$ and a set of estimated connection strength $b_{ij}$.

In cases of no latent confounders, Algorithm 3 is essentially equivalent to DirectLiNGAM (Shimizu et al., 2011). Matlab codes for performing Algorithm 3 are available at `http://www.ar.sanken.osaka-u.ac.jp/~sshimizu/code/Plingamcode.html`.

# 4 Experiments on artificial data

We compared our method with two estimation methods for LiNGAM in Eq. (2) called ICA-LiNGAM (Shimizu et al., 2006) and DirectLiNGAM (Shimizu et al., 2011) that do not allow latent confounders and an estimation method for latent variable LiNGAM in Eq. (3) called Pairwise LvLiNGAM (Entner & Hoyer, 2011). If there are no latent confounders, all the methods should estimate correct causal orders for large enough sample sizes. The numbers of variables were 5, 10, and 15, and the sample sizes tested were 500, 1000, and 1500. The original networks used were shown in Fig. 2 to Fig. 4. The $e_1$, $e_4$, $e_7$, $e_{10}$, $e_{13}$, $f_1$ and $f_4$ followed a multimodal asymmetric mixture of two Gaussians, $e_2$, $e_5$, $e_8$, $e_{11}$, $e_{14}$, $f_2$ and $f_5$ followed a double exponential distribution, and $e_3$, $e_6$, $e_9$, $e_{12}$, $e_{15}$, $f_3$ and $f_6$ followed a multimodal symmetric mixture of two Gaussians. The variances of the $e_i$ were set so that $\mathrm{var}(e_i)/\mathrm{var}(x_i){=}1/2$. We permuted the variables according to a random ordering. The number of trials was 100. The significance level $\alpha$ was 0.05.

First, to evaluate performance of estimating causal orders $k(i)$, we computed the percentage of correctly estimated causal orders among estimated causal orders between two variables (Precision) and the percentage of correctly estimated causal orders among actual causal orders between two variables (Recall). We also computed the F-measure defined as $2 \times$ Precision $\times$ Recall/(Precision + Recall), which is the harmonic mean of Precision and Recall. The reason why only pairwise causal orders were evaluated was that Pairwise LvLiNGAM only estimates causal orders of two variables unlike our method and DirectLiNGAM. Tables 1, 2 and 3 show the results. Regarding recalls and F-measures, the maximal performances when no statistical errors occur are also shown in the right-most columns. For example in Fig. 2, Pairwise LvLiNGAM can find all the causal orderings except $k(2) < k(4)$, $k(2) < k(5)$, $k(3) < k(4)$ and $k(3) < k(5)$. ParceLiNGAM further can find $k(2) < k(5)$ and $k(3) < k(5)$ since it estimates causal orderings between more than two variables. In some cases, the empirical recalls and F-measures were higher than their maximal performances. This is because causal orders of some variables that are affected by latent confounders happened to be correctly estimated. Regarding precisions and F-measures, our method ParceLiNGAM worked best for all the conditions. Regarding recalls,
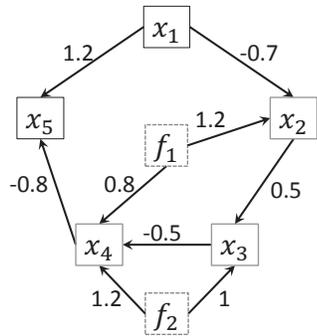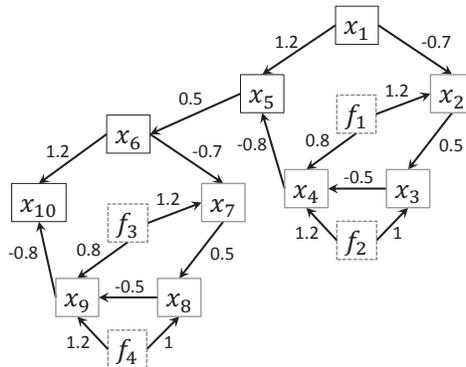
Figure 2: 5 variable network
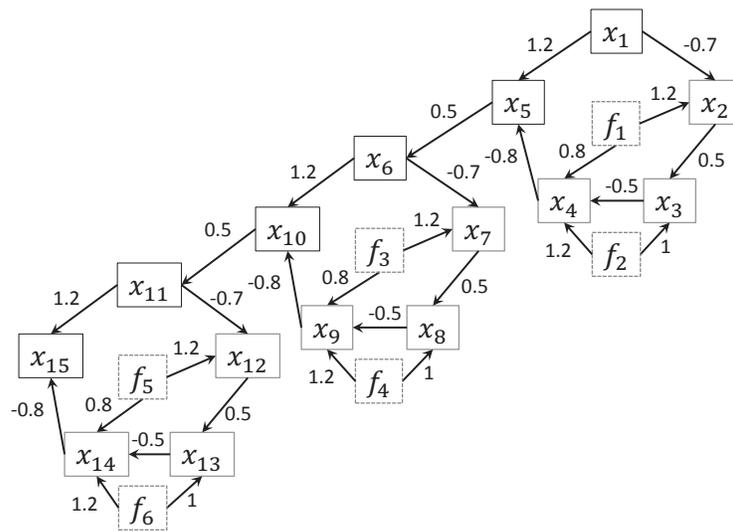


Figure 3: 10 variable network



Figure 4: 15 variable network

ParceLiNGAM worked best for most conditions and was the second-best but comparable to the best method DirectLiNGAM for the other conditions.

Next, to evaluate the performance in estimating connection strengths $b_{ij}$, we computed the root mean square errors between true connection strengths and estimated ones. Note that Pairwise LvLiNGAM does not estimate $b_{ij}$. Table 4 show the results. Our method was most accurate for all the conditions.

Table 5 shows average computation times. The amount out computation of our ParceLiNGAM was larger than the other methods when the sample size was increased. However, its amount of computation can be considered to be still tractable. For larger numbers of variables, we would need to select a subset of variables to decrease the number of variables to be analyzed. However, this selection does not bias results of our method since it allows latent confounders.

Table 1: Precisions

|  |  | Sample size | | |
|---|---|---|---|---|
|  |  | 500 | 1000 | 1500 |
| ParceLiNGAM | dim.=5 | 1.0 | 1.0 | 1.0 |
|  | dim.=10 | 0.81 | 0.88 | 0.93 |
|  | dim.=15 | 0.81 | 0.89 | 0.92 |
| PairwiseLvLiNGAM | dim.=5 | 0.87 | 0.94 | 0.94 |
|  | dim.=10 | 0.75 | 0.79 | 0.81 |
|  | dim.=15 | 0.67 | 0.76 | 0.75 |
| DirectLiNGAM | dim.=5 | 0.82 | 0.88 | 0.85 |
|  | dim.=10 | 0.59 | 0.71 | 0.73 |
|  | dim.=15 | 0.78 | 0.80 | 0.82 |
| ICA-LiNGAM | dim.=5 | 0.80 | 0.75 | 0.76 |
|  | dim.=10 | 0.62 | 0.62 | 0.58 |
|  | dim.=15 | 0.58 | 0.59 | 0.58 |

Table 2: Recalls

|  |  | Sample size | | | Max. performance |
|---|---|---|---|---|---|
|  |  | 500 | 1000 | 1500 |  |
| ParceLiNGAM | dim.=5 | 0.86 | 0.82 | 0.80 | 0.80(8/10) |
|  | dim.=10 | 0.79 | 0.85 | 0.91 | 0.91(41/45) |
|  | dim.=15 | 0.80 | 0.87 | 0.89 | 0.94(99/105) |
| PairwiseLvLiNGAM | dim.=5 | 0.65 | 0.62 | 0.59 | 0.60(6/10) |
|  | dim.=10 | 0.50 | 0.55 | 0.54 | 0.49(22/45) |
|  | dim.=15 | 0.39 | 0.45 | 0.43 | 0.46(48/105) |
| DirectLiNGAM | dim.=5 | 0.82 | 0.88 | 0.85 | - |
|  | dim.=10 | 0.59 | 0.71 | 0.73 | - |
|  | dim.=15 | 0.78 | 0.80 | 0.82 | - |
| ICA-LiNGAM | dim.=5 | 0.80 | 0.75 | 0.76 | - |
|  | dim.=10 | 0.62 | 0.62 | 0.58 | - |
|  | dim.=15 | 0.58 | 0.59 | 0.58 | - |

# 5 Experiments on simulated fMRI data

Finally, we tested our method on simulated functional magnetic resonance imaging (fMRI) data generated in (Smith et al., 2011) based on a well-known mathematical brain model called the dynamic causal modeling (Friston, Harrison, & Penny,

Table 3: F-measures

| | | Sample size | | | Max. performance |
|---|---|---|---|---|---|
| | | 500 | 1000 | 1500 | |
| ParceLiNGAM | dim.=5 | 0.92 | 0.90 | 0.89 | 0.89 |
| | dim.=10 | 0.80 | 0.86 | 0.92 | 0.95 |
| | dim.=15 | 0.81 | 0.88 | 0.90 | 0.97 |
| PairwiseLvLiNGAM | dim.=5 | 0.75 | 0.75 | 0.72 | 0.75 |
| | dim.=10 | 0.60 | 0.65 | 0.65 | 0.66 |
| | dim.=15 | 0.49 | 0.56 | 0.54 | 0.63 |
| DirectLiNGAM | dim.=5 | 0.82 | 0.88 | 0.85 | - |
| | dim.=10 | 0.59 | 0.71 | 0.73 | - |
| | dim.=15 | 0.78 | 0.80 | 0.82 | - |
| ICA-LiNGAM | dim.=5 | 0.80 | 0.75 | 0.76 | - |
| | dim.=10 | 0.62 | 0.62 | 0.58 | - |
| | dim.=15 | 0.58 | 0.59 | 0.58 | - |

Table 4: Root Mean Square Errors

| | | Sample size | | |
|---|---|---|---|---|
| | | 500 | 1000 | 1500 |
| ParceLiNGAM | dim.=5 | 0.030 | 0.020 | 0.016 |
| | dim.=10 | 0.078 | 0.060 | 0.052 |
| | dim.=15 | 0.083 | 0.046 | 0.031 |
| DirectLiNGAM | dim.=5 | 0.22 | 0.16 | 0.18 |
| | dim.=10 | 0.16 | 0.083 | 0.089 |
| | dim.=15 | 0.096 | 0.074 | 0.070 |
| ICA-LiNGAM | dim.=5 | 0.11 | 0.11 | 0.10 |
| | dim.=10 | 0.16 | 0.15 | 0.15 |
| | dim.=15 | 0.16 | 0.14 | 0.13 |

Table 5: Computational Times

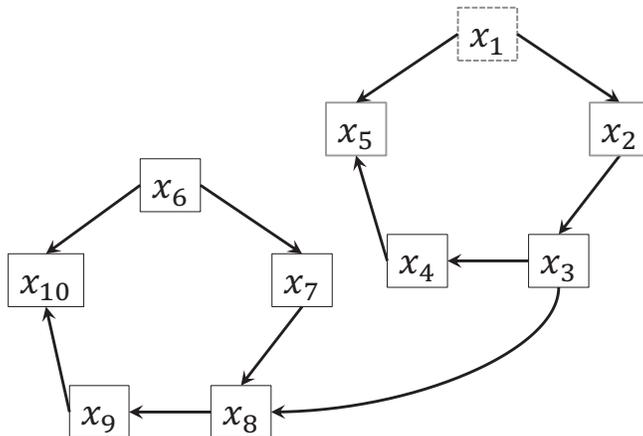| | | Sample size | | |
|---|---|---|---|---|
| | | 500 | 1000 | 1500 |
| ParceLiNGAM | dim.=5 | 0.66 sec. | 1.7 sec. | 4.4 sec. |
| | dim.=10 | 10 sec. | 1.5 min. | 8.1 min. |
| | dim.=15 | 8.5 min. | 5.3 hrs. | 19 hrs. |
| PairwiseLvLiNGAM | dim.=5 | 0.64 sec. | 2.6 sec. | 7.0 sec. |
| | dim.=10 | 2.8 sec. | 12 sec. | 30 sec. |
| | dim.=15 | 6.6 sec. | 29 sec. | 74 sec. |
| DirectLiNGAM | dim.=5 | 0.23 sec. | 0.84 sec. | 1.2 sec. |
| | dim.=10 | 1.7 sec. | 7.3 sec. | 11 sec. |
| | dim.=15 | 6.4 sec. | 29 sec. | 44 sec. |
| ICA-LiNGAM | dim.=5 | 0.12 sec. | 0.051 sec. | 0.047 sec. |
| | dim.=10 | 0.34 sec. | 0.18 sec. | 0.10 sec. |
| | dim.=15 | 0.81 sec. | 0.68 sec. | 0.53 sec. |

Figure 5: The network used in the simulated fMRI experiments. We omitted $x_1$ to create a latent confounder.

2003). We used Simulation 2 data and Simulation 6 data. Both datasets consisted of 10 variables whose causal structure is shown in Fig. 5. The session durations were 10 minutes (200 time points) and 60 minutes (1200 time points), respectively. We also created a dataset of 30 minutes (600 time points) by taking the first half of Simulation 6 data. Although these data are time-series, we did not add lag-based approaches including vector autoregressive models into comparison as in (Hyvärinen & Smith, 2013) since it was shown by (Smith et al., 2011) that lag-based methods worked poorly on these Simulation 2 data and Simulation 6 data.

For each of the three different duration settings, we gave the 50 datasets (one by one) to ParceLiNGAM, PairwiseLvLiNGAM, DirectLiNGAM and ICA-LiNGAM after omitting $x_1$ to create a latent confounder and randomly permuting the other variables. Table 6 shows the precision, recalls, and F-measures of causal orders. Regarding precisions, we excluded such variable pairs $x_i$ and $x_j$ that one has no directed path to the other, e.g., $x_2$ and $x_6$, since both $k(i) < k(j)$ and $k(i) > k(j)$ are correct. This was because estimation of causal directions is the main topic of this paper. The significance level $\alpha$ was 0.05. For all of the cases, ParceLiNGAM worked better than the others.

## 6    Conclusions

We proposed a new algorithm for learning causal orders, which is robust against latent confounders. In experiments on artificial data and simulated fMRI data,

Table 6: Results on simulated fMRI data

|  |  | sim2 (10 min.) | sim6 (30 min.) | sim6 (60 min.) |
|---|---|---|---|---|
| ParceLiNGAM | Precision | 0.54 | 0.56 | 0.60 |
|  | Recall | 0.53 | 0.55 | 0.58 |
|  | F-measure | 0.53 | 0.55 | 0.59 |
| PairwiseLvLiNGAM | Precision | 0.31 | 0.25 | 0.24 |
|  | Recall | 0.22 | 0.15 | 0.14 |
|  | F-measure | 0.26 | 0.19 | 0.18 |
| DirectLiNGAM | Precision | 0.50 | 0.51 | 0.45 |
|  | Recall | 0.50 | 0.51 | 0.45 |
|  | F-measure | 0.50 | 0.51 | 0.45 |
| ICA-LiNGAM | Precision | 0.49 | 0.47 | 0.47 |
|  | Recall | 0.49 | 0.47 | 0.47 |
|  | F-measure | 0.49 | 0.47 | 0.47 |

our methods learned more causal orders correctly than existing methods. An important problem for future research is to develop computationally more efficient algorithms. One approach might be to develop a divide-and-conquer algorithm that divides variables into subsets with moderate numbers of variables and integrates the estimation results on the subsets.

### Acknowledgments.

# References

Bollen, K. (1989). *Structural equations with latent variables*. John Wiley & Sons.

Chen, Z., & Chan, L. (2013). Causality in linear nongaussian acyclic models in the presence of latent gaussian confounders. *Neural Computation*. (In press.)

Darmois, G. (1953). Analyse générale des liaisons stochastiques. *Review of the International Statistical Institute*, *21*, 2-8.

Entner, D., & Hoyer, P. O. (2011). Discovering unconfounded causal relationships using linear non-gaussian models. In *New Frontiers in Artificial Intelligence, Lecture Notes in Computer Science* (Vol. 6797, p. 181-195).

Fisher, R. (1950). *Statistical methods for research workers*. Oliver and Boyd.

Friston, K., Harrison, L., & Penny, W. (2003). Dynamic causal modelling. *Neuroimage*, *19*(4), 1273–1302.

Gretton, A., Fukumizu, K., Teo, C., Song, L., Schölkopf, B., & Smola, A. J. (2008). A kernel statistical test of independence. In *Advances in Neural Information Processing Systems 20*. Cambridge, MA: MIT Press.

Hoyer, P. O., Janzing, D., Mooij, J., Peters, J., & Schölkopf, B. (2009). Nonlinear causal discovery with additive noise models. In *Advances in Neural Information Processing Systems 21* (pp. 689–696).

Hoyer, P. O., Shimizu, S., Kerminen, A., & Palviainen, M. (2008). Estimation of causal effects using linear non-gaussian causal models with hidden variables. *International Journal of Approximate Reasoning*, *49*(2), 362-378.

Hyvärinen, A., Karhunen, J., & Oja, E. (2001). *Independent component analysis.* New York: Wiley.

Hyvärinen, A., & Smith, S. M. (2013). Pairwise likelihood ratios for estimation of non-gaussian structural equation models. *Journal of Machine Learning Research*, *14*, 111–152.

Lewicki., M., & Sejnowski, T. J. (2000). Learning overcomplete representations. *Neural Computation*, *12*(2), 337-365.

Mooij, J., Janzing, D., Peters, J., & Schölkopf, B. (2009). Regression by dependence minimization and its application to causal inference in additive noise models. In *Proc. the 26th Int. Conf. on Machine Learning (ICML2009)* (p. 745-752).

Pearl, J. (2000). *Causality: Models, reasoning, and inference.* Cambridge University Press. ((2nd ed. 2009))

Shimizu, S., Hoyer, P. O., Hyvärinen, A., & Kerminen, A. (2006). A linear non-gaussian acyclic model for causal discovery. *J. Mach. Learn. Res.*, *7*, 2003-2030.

Shimizu, S., Inazumi, T., Sogawa, Y., Hyvärinen, A., Kawahara, Y., Washio, T., et al. (2011). DirectLiNGAM: A direct method for learning a linear non-Gaussian structural equation model. *J. Mach. Learn. Res.*, *12*, 1225-1248.

Skitovitch, W. P. (1953). On a property of the normal distribution. *Doklady Akademii Nauk SSSR*, *89*, 217-219.

Smith, S. M., Miller, K. L., Salimi-Khorshidi, G., Webster, M., Beckmann, C. F., Nichols, T. E., et al. (2011). Network modelling methods for FMRI. *NeuroImage*, *54*(2), 875–891.

Spirtes, P., Glymour, C., & Scheines, R. (1993). *Causation, prediction, and search.* Springer Verlag. ((2nd ed. MIT Press 2000))

Tashiro, T., Shimizu, S., Hyvärinen, A., & Washio, T. (2012). Estimation of causal orders in a linear non-gaussian acyclic model: a method robust against latent confounders. In *Proc. Int. Conf. on Artificial Neural Networks (ICANN2012), Lausanne, Switzerland* (p. 491-498).

# Appendix

We first give Darmois-Skitovitch theorem (Darmois, 1953; Skitovitch, 1953):

**Theorem 1 (Darmois-Skitovitch theorem (D-S theorem))** *Define two random variables $y_1$ and $y_2$ as linear combinations of independent random variables $s_i(i=1, \cdots, q)$: $y_1 = \sum_{i=1}^{q} \alpha_i s_i$, $y_2 = \sum_{i=1}^{q} \beta_i s_i$. Then, if $y_1$ and $y_2$ are independent, all variables $s_j$ for which $\alpha_j \beta_j \neq 0$ are Gaussian.* □

In other words, this theorem means that if there exists a non-Gaussian $s_j$ for which $\alpha_j \beta_j \neq 0$, $y_1$ and $y_2$ are dependent.

**Proof of Lemma 1**

i) Assume that $x_j$ has at least one parent observed variable or latent confounder. Let $P_j$ denote the set of the parent variables of $x_j$. Then one can write $x_j = \sum_{p_h \in P_j} w_{jh} p_h + e_j$, where the parent variables $p_h$ are independent of $e_j$ and the coefficients $w_{jh}$ are non-zero. Let a vector $\mathbf{x}_{P_j}$ and a column vector $\boldsymbol{w}_{P_j}$ collect all the variables in $P_j$ and the corresponding connection strengths, respectively. Then, the covariances between $\mathbf{x}_{P_j}$ and $x_j$ are $E(\mathbf{x}_{P_j} x_j) = E\{\mathbf{x}_{P_j}(\boldsymbol{w}_{P_j}^T \mathbf{x}_{P_j} + e_j)\} = E(\mathbf{x}_{P_j} \mathbf{x}_{P_j}^T)\boldsymbol{w}_{P_j}$. The covariance matrix $E(\mathbf{x}_{P_j} \mathbf{x}_{P_j}^T)$ is positive definite since the external influences and latent confounders are mutually independent and have positive variances. Thus, the covariance vector $E(\mathbf{x}_{P_j} x_j) = E(\mathbf{x}_{P_j} \mathbf{x}_{P_j}^T)\boldsymbol{w}_{P_j}$ above cannot equal the zero vector, and there must be at least one variable in $P_j$ with which $x_j$ covaries.

i-a) Suppose that $x_i$ is a parent of $x_j$ in $P_j$ that covaries with $x_j$. For such $x_i$, we have

$$
\begin{align}
r_i^{(j)} &= x_i - \frac{\mathrm{cov}(x_i, x_j)}{\mathrm{var}(x_j)} x_j \tag{8} \\
&= x_i - \frac{\mathrm{cov}(x_i, x_j)}{\mathrm{var}(x_j)} \Big( \sum_{p_h \in P_j} w_{jh} p_h + e_j \Big) \tag{9} \\
&= \Big\{ 1 - \frac{w_{ji}\mathrm{cov}(x_i, x_j)}{\mathrm{var}(x_j)} \Big\} x_i - \frac{\mathrm{cov}(x_i, x_j)}{\mathrm{var}(x_j)} \sum_{p_h \in P_j, p_h \neq x_i} w_{jh} p_h \\
&\quad - \frac{\mathrm{cov}(x_i, x_j)}{\mathrm{var}(x_j)} e_j. \tag{10}
\end{align}
$$

Each of those parent variables (including $x_i$) in $P_j$ is a linear combination of external influences *other than* $e_j$ and latent confounders that are non-Gaussian and independent. Thus, the $r_i^{(j)}$ and $x_j$ can be written as linear combinations of non-Gaussian and independent external influences including $e_j$ and latent confounders. Further, the coefficient of $e_j$ on $r_i^{(j)}$ is non-zero since $\mathrm{cov}(x_i, x_j) \neq 0$ aforementioned and that on $x_j$ is one by definition. These imply that $r_i^{(j)}$ and $x_j$ are dependent since $r_i^{(j)}$, $x_j$ and $e_j$ correspond to $y_1$, $y_2$, $s_j$ in D-S theorem, respectively.

i-b) Next, suppose that $x_j$ has a latent confounder $f_k$ in $P_j$ that covaries with $x_j$. The latent confounder $f_k$ should have a non-zero coefficient on at least one other observed variable $x_i$. Without loss of generality, it is enough to consider two observed variable cases that we only observe $x_i$ and $x_j$:

$$
\begin{align}
x_i &= b_{ij} x_j + \lambda_{ik} f_k + e_i + \sum_{h \neq k} \lambda_{ih} f_h \tag{11} \\
x_j &= b_{ji} x_i + \lambda_{jk} f_k + e_j + \sum_{l \neq k} \lambda_{il} f_l, \tag{12}
\end{align}
$$

where $\lambda_{ik}$ and $\lambda_{jk}$ are non-zero since $f_k$ is a latent confounder of $x_i$ and $x_j$. Since the model is acyclic, $b_{ij} b_{ji} = 0$.

16

First, suppose that $b_{ij}$ is zero. Then, we have

$$
\begin{aligned}
r_i^{(j)} &= x_i - \frac{\text{cov}(x_i, x_j)}{\text{var}(x_j)} x_j & (13)\\
&= \{\lambda_{ik} - \frac{\text{cov}(x_i, x_j)}{\text{var}(x_j)}(b_{ji}\lambda_{ik} + \lambda_{jk})\}f_k \\
&\quad +\{1 - \frac{\text{cov}(x_i, x_j)}{\text{var}(x_j)}b_{ji}\}e_i - \frac{\text{cov}(x_i, x_j)}{\text{var}(x_j)}e_j + D_1, & (14)
\end{aligned}
$$

where $D_1$ is a linear combinations of non-Gaussian and independent latent confounders other than $f_k$. If $\text{cov}(x_i, x_j)$ is zero, the coefficient of $f_k$ on $r_i^{(j)}$ is $\lambda_{ik}$ and is non-zero. If $\text{cov}(x_i, x_j)$ is non-zero, the coefficient of $e_j$ on $r_i^{(j)}$ is $-\frac{\text{cov}(x_i, x_j)}{\text{var}(x_j)}$ and is non-zero. Thus, in both of the cases, $r_i^{(j)}$ and $x_j$ are dependent due to D-S theorem. Remember that the coefficient of $e_j$ on $x_j$ is one by definition.

Next, suppose that $b_{ji}$ is zero. Then, we have

$$
\begin{aligned}
r_i^{(j)} &= x_i - \frac{\text{cov}(x_i, x_j)}{\text{var}(x_j)} x_j & (15)\\
&= \{(b_{ij}\lambda_{jk} + \lambda_{ik}) - \frac{\text{cov}(x_i, x_j)}{\text{var}(x_j)}\lambda_{jk}\}f_k \\
&\quad + e_i + (b_{ij} - \frac{\text{cov}(x_i, x_j)}{\text{var}(x_j)})e_j + D_2, & (16)
\end{aligned}
$$

where $D_2$ is a linear combinations of non-Gaussian and independent latent confounders other than $f_k$. If $\text{cov}(x_i, x_j)$ is zero and $b_{ij}$ is zero, the coefficient of $f_k$ on $r_i^{(j)}$ is $\lambda_{ik}$ and is non-zero. If $\text{cov}(x_i, x_j)$ is zero and $b_{ij}$ is non-zero, the coefficient of $e_j$ on $r_i^{(j)}$ is $b_{ij}$ and is non-zero. If $\text{cov}(x_i, x_j)$ is non-zero and $b_{ij}$ is zero, the coefficient of $e_j$ on $r_i^{(j)}$ is $-\frac{\text{cov}(x_i, x_j)}{\text{var}(x_j)}$ and is non-zero. If $\text{cov}(x_i, x_j)$ is non-zero and $b_{ij}$ is non-zero, either of the followings holds: a) the coefficient of $e_j$ on $r_i^{(j)}$ is non-zero, that is, $b_{ij} \neq \text{cov}(x_i, x_j)/\text{var}(x_j)$ or b) the coefficient of $e_j$ on $r_i^{(j)}$ is zero and hence the coefficient of $f_k$ on $r_i^{(j)}$ is $\lambda_{ik}$ and is non-zero. Thus, in all of the cases, $r_i^{(j)}$ and $x_j$ are dependent due to D-S theorem.

ii) The converse of contrapositive of i) is straightforward using the model definition. From i) and ii), the lemma is proven. ∎

**Proof of Lemma 2**

i) Assume that a variable $x_j$ has at least one child observed variable or latent confounder. First, without loss of generality, one can write

$$
\begin{aligned}
\boldsymbol{x} = \begin{bmatrix} x_j \\ \boldsymbol{x}_{(-j)} \end{bmatrix} &= (\mathbf{I} - \mathbf{B})^{-1}(\boldsymbol{\Lambda}\boldsymbol{f} + \boldsymbol{e}) = \mathbf{A}(\boldsymbol{\Lambda}\boldsymbol{f} + \boldsymbol{e}) & (17)\\
&= \begin{bmatrix} 1 & \boldsymbol{a}_{j(-j)}^T \\ \boldsymbol{a}_{(-j)j} & \mathbf{A}_{(-j)} \end{bmatrix} \begin{bmatrix} \boldsymbol{\lambda}_j^T \boldsymbol{f} + e_j \\ \boldsymbol{\Lambda}_{(-j)}\boldsymbol{f} + \boldsymbol{e}_{(-j)} \end{bmatrix}, & (18)
\end{aligned}
$$

where each of $\mathbf{A}$ $(= (\mathbf{I} - \mathbf{B})^{-1})$ and $\mathbf{A}_{(-j)}$ is invertible and can be permuted to be a lower triangular matrix with the diagonal elements being ones if the rows and columns are simultaneously permuted according to the causal ordering $k(i)$. The same applies to the inverse of $\mathbf{A}$:

$$\mathbf{A}^{-1} = \begin{bmatrix} (1 - \boldsymbol{a}_{j(-j)}^T \mathbf{A}_{(-j)}^{-1} \boldsymbol{a}_{(-j)j})^{-1} & -\boldsymbol{a}_{j(-j)}^T \mathbf{D}^{-1} \\ -\mathbf{D}^{-1} \boldsymbol{a}_{(-j)j} & \mathbf{D}^{-1} \end{bmatrix}, \tag{19}$$

where $\mathbf{D} = \mathbf{A}_{(-j)} - \boldsymbol{a}_{(-j)j} \boldsymbol{a}_{j(-j)}^T$. Thus, $1 - \boldsymbol{a}_{j(-j)}^T \mathbf{A}_{(-j)}^{-1} \boldsymbol{a}_{(-j)j} = 1$.

Then,

$$
\begin{aligned}
r_j^{(-j)} &= x_j - \boldsymbol{\sigma}_{(-j)j}^T \Sigma_{(-j)}^{-1} \boldsymbol{x}_{(-j)} \tag{20} \\
&= \boldsymbol{\lambda}_j^T \boldsymbol{f} + e_j + \boldsymbol{a}_{j(-j)}^T (\boldsymbol{\Lambda}_{(-j)} \boldsymbol{f} + \boldsymbol{e}_{(-j)}) \\
&\quad - \boldsymbol{\sigma}_{(-j)j}^T \Sigma_{(-j)}^{-1} \{ \boldsymbol{a}_{(-j)j} (\boldsymbol{\lambda}_j^T \boldsymbol{f} + e_j) + \mathbf{A}_{(-j)} (\boldsymbol{\Lambda}_{(-j)} \boldsymbol{f} + \boldsymbol{e}_{(-j)}) \tag{21} \\
&= \{ \boldsymbol{\lambda}_j^T + \boldsymbol{a}_{j(-j)}^T \boldsymbol{\Lambda}_{(-j)} - \boldsymbol{\sigma}_{(-j)j}^T \Sigma_{(-j)}^{-1} (\boldsymbol{a}_{(-j)j} \boldsymbol{\lambda}_j^T + \mathbf{A}_{(-j)} \boldsymbol{\Lambda}_{(-j)}) \} \boldsymbol{f} \\
&\quad + \{ 1 - \boldsymbol{\sigma}_{(-j)j}^T \Sigma_{(-j)}^{-1} \boldsymbol{a}_{(-j)j} \} e_j + \{ \boldsymbol{a}_{j(-j)}^T - \boldsymbol{\sigma}_{(-j)j}^T \Sigma_{(-j)}^{-1} \mathbf{A}_{(-j)} \} \boldsymbol{e}_{(-j)} \tag{22}
\end{aligned}
$$

In Eq.(22), if $\boldsymbol{a}_{j(-j)}^T - \boldsymbol{\sigma}_{(-j)j}^T \Sigma_{(-j)}^{-1} \mathbf{A}_{(-j)} = \boldsymbol{0}^T$, then we have

$$
\begin{aligned}
r_j^{(-j)} &= \{ \boldsymbol{\lambda}_j^T (1 - \boldsymbol{a}_{j(-j)}^T \mathbf{A}_{(-j)}^{-1} \boldsymbol{a}_{(-j)j}) \} \boldsymbol{f} + \{ 1 - \boldsymbol{a}_{j(-j)}^T \mathbf{A}_{(-j)}^{-1} \boldsymbol{a}_{(-j)j} \} e_j \tag{23} \\
&= \boldsymbol{\lambda}_j^T \boldsymbol{f} + e_j. \tag{24}
\end{aligned}
$$

Thus, the coefficient of $e_j$ on $r_j^{(-j)}$ is one. Now, suppose that $x_j$ has a child $x_i$. If the coefficient of $e_j$ on $x_i$ is non-zero, $r_j^{(-j)}$ and $\boldsymbol{x}_{(-j)}$ are dependent due to D-S theorem. Even if it is zero, $i.e.$, cancelled out to be zero by special parameter values of the connection strengths, the coefficient of $e_j$ on at least one other variable in $\boldsymbol{x}_{(-j)}$ is non-zero since there must be such an observed variable to cancel out the coefficient of $e_j$ on $x_i$ to be zero. It implies that $r_j^{(-j)}$ and $\boldsymbol{x}_{(-j)}$ are dependent due to D-S theorem. Next, suppose that $x_j$ has a latent confounder $f_i$. Then, in Eq.(24), the corresponding element in $\boldsymbol{\lambda}_j$ is not zero, $i.e.$, the coefficient of $f_i$ on $r_j^{(-j)}$ is not zero. Further, $f_i$ has a non-zero coefficient on at least one variable in $\boldsymbol{x}_{(-j)}$ due to the definition of latent confounders. Therefore, $r_j^{(-j)}$ and $\boldsymbol{x}_{(-j)}$ are dependent due to D-S theorem.

On the other hand, in Eq.(22), if $\boldsymbol{a}_{j(-j)}^T - \boldsymbol{\sigma}_{(-j)j}^T \Sigma_{(-j)}^{-1} \mathbf{A}_{(-j)} \neq \boldsymbol{0}^T$, at least one of the coefficients of the elements in $\boldsymbol{e}_{(-j)}$ on $r_j^{(-j)}$ is not zero. By definition, every element in $\boldsymbol{e}_{(-j)}$ has a non-zero coefficient on the corresponding element in $\boldsymbol{x}_{(-j)}$, Thus, $r_j^{(-j)}$ and $\boldsymbol{x}_{(-j)}$ are dependent due to D-S theorem.

ii) The converse of contrapositive of i) is straightforward using the model definition. From i) and ii), the lemma is proven. ∎