

# Technical Proof for ANOCE: Analysis of Causal Effects with Multiple Mediators via Constrained Structural Learning

## 1 Connection to Literature

We establish the connection between our proposed method to the literature from three different angles. First, we show that the individual mediation effect defined in Chakraborty et al. [2] can be decomposed into our defined  $DM$  and  $IM$  when the LSEM assumption holds. Next, we give an equivalent definition of the  $DM$  through a type of special edge (last edge) in the causal graph. Lastly, we prove that the proposed  $DM$  is consistent with the interventional effect via a particular mediator defined in Vansteelandt and Daniel [5] under the LSEM.

### 1.1 From Individual Mediation Viewpoint

Chakraborty et al. [2] defined the individual mediation effect under the LSEM as follows.

**Definition 1.1** [2] *Individual mediation effect for  $M_i$ :*

$$\eta_i = \left[ E\{M_i | do(A = a + 1)\} - E\{M_i | do(A = a)\} \right] \times \left[ E\{Y | do(M_i = m_i + 1)\} - E\{Y | do(M_i = m_i)\} \right]. \quad (1)$$

In the following theorem, we show that the summation of  $\eta_i$  is strictly larger than the  $IE$  if the mediators are not parallel. The proof is given in Section 2.2.

**Theorem 1.1** *If there exists at least one directed path  $\pi^* \in \{\pi_{AY}(\mathcal{G})\}$  such that the length of  $\pi^*$  is larger than 2, and the element in  $B$  is nonnegative, then*

$$\sum \eta_i > IE. \quad (2)$$

**Remark 1.1** *From the above theorem, it is clear that the mediator effect defined in Chakraborty et al. [2] is not appropriate for interpreting the decomposition of the indirect effect, when there exists interaction among mediators (a common situation as described in the introduction). Here, we keep the condition that the element in  $B$  is nonnegative, as the multiple count mediation effects in Chakraborty et al. [2] may cancel out in some cases and their summation would equal to the  $IE$  by chance.*

Inspired by the proof of Theorem 1.1, the mediator effect  $\eta_i$  can be decomposed into two parts, the natural direct and indirect effect for  $i$ -th mediator, as shown in the following corollary.

**Corollary 1.1** *Under assumptions (A1-A3) and Model 1, we have*

$$\eta_i = DM_i + IM_i. \quad (3)$$

**Remark 1.2** *Corollary 1.1 together with the definition of  $\eta_i$  in Chakraborty et al. [2] provides a feasible way to numerically calculate the natural indirect effect  $IM_i$ . Specifically, by deleting the mediator  $M_i$  in the causal graph, the reduced treatment effect corresponds to  $\eta_i$ , then  $IM_i = \eta_i - DM_i$ , where the explicit expression of the  $DM_i$  is provided in Theorem 3.1. See more implementation details in Section 4.*

## 1.2 From Graphical Perspective

Next, we give the definition of the edge-specific effect following Avin et al. [1]. Suppose a directed edge of interest as  $X_i \rightarrow X_j$  in a weighted DAG  $\mathcal{G}$ . Define a new weighted DAG  $\mathcal{G}'_{i,j}$  by deleting the directed edge  $X_i \rightarrow X_j$  in  $\mathcal{G}$ , i.e.,  $\mathcal{G}'_{i,j} \equiv \mathcal{G} \setminus (X_i \rightarrow X_j)$ .

**Definition 1.2** [1] *Edge-specific effect:*

$$ET(X_i, X_j) = TE_{\mathcal{G}} - TE_{\mathcal{G}'_{i,j}}, \quad (4)$$

where  $TE_{\mathcal{G}}$  means the total effect in graph  $\mathcal{G}$ .

We next give an equivalent definition of our proposed  $DM$  from a graphical perspective. Let the edge in  $\mathcal{G}$  that starts with  $i$ -th mediator and ends with node  $Y$ , i.e.,  $M_i \rightarrow Y$ , as the  $i$ -th last edge. Denote the graph  $\mathcal{G}$  deleting the  $i$ -th last edge ( $M_i \rightarrow Y$ ) as  $\mathcal{G}'_i$ . We define the  $i$ th last edge-specific effect as

**Definition 1.3** *Last edge-specific effect for  $M_i$ :*

$$LE_i = \begin{cases} TE_{\mathcal{G}} - TE_{\mathcal{G}'_i}, & \text{if there exists edge } M_i \rightarrow Y \text{ in } \mathcal{G}; \\ 0, & \text{otherwise.} \end{cases} \quad (5)$$

By Theorem 3.1, we have  $(I_p - B_M^\top)^{-1}\alpha$  is the causal effect of  $A$  on  $M$ . Let  $\zeta \equiv (I_p - B_M^\top)^{-1}\alpha$ , with its  $i$ -th element  $\zeta_i \equiv \{(I_p - B_M^\top)^{-1}\alpha\}_i$ . Next, we show that the  $i$ -th last edge-specific effect can be presented as  $\beta_i\zeta_i$  under the LSEM in the following theorem, where  $\beta_i$  is the  $i$ -th element of the vector  $\beta$  and corresponds to the weight of the edge  $M_i \rightarrow Y$ . The proof can be found in Section 2.3.

**Theorem 1.2** *Under assumptions (A1-A3) and Model 1, we have*

$$LE_i = \beta_i\zeta_i. \quad (6)$$

Based on Theorem 3.1, the natural direct effect of  $M_i$  on  $Y$  can be expressed as  $DM_i = \beta_i\zeta_i$ . Thus, with the result of Theorem 1.2, it is easy to show the following corollary.

**Corollary 1.2** *Under assumptions (A1-A3) and Model 1, the natural direct effect of  $M_i$  is equal to the  $i$ -th last edge-specific effect:*

$$LE_i = DM_i = \beta_i \zeta_i. \quad (7)$$

**Remark 1.3** *Here, both definitions describe the direct impact of one mediator  $M_i$  on the outcome. The natural direct effect of a particular mediator  $M_i$  can be understood as the influence when removing the direct edge between  $M_i$  and  $Y$ . Thus, we have the equivalence between two definitions.*

Then, we can decompose the total natural indirect effect into  $p$  last edge-specific effects or  $p$  DMs as the following additive form, based on Theorem 3.2 and Corollary 1.2.

**Corollary 1.3** *Under assumptions (A1-A3) and Model 1, we have*

$$IE = \beta^\top \zeta = \sum_{i=1}^p \beta_i \zeta_i = \sum_{i=1}^p DM_i = \sum_{i=1}^p LE_i. \quad (8)$$

In fact, based on the uniqueness of each last edge, the natural indirect effect can be decomposed into  $p$  last edge-specific effect regardless of the LSEM setting through the graphical perspective. We give the following intuitive conclusion. The proof can be found in Section 2.4.

**Theorem 1.3** *The IE can be decomposed through LEs as:*

$$IE = \sum_{i=1}^p LE_i. \quad (9)$$

**Remark 1.4** *One can view the last edge-specific effect as the generalized definition of the natural direct effect for mediator without the LSEM assumption.*

### 1.3 From Interventional Effect Level

Finally, we show the consistency of our defined  $DM$  to the interventional effect via a particular mediator defined in Vansteelandt and Daniel [5] under the LSEM.

**Definition 1.4** [5] *Under assumptions (A1-A3), the interventional effect via  $M_i$  is*

$$\begin{aligned} \xi_i = \sum_{m_1 \in \mathbb{M}_1} \cdots \sum_{m_p \in \mathbb{M}_p} & \left[ E(Y|A = a, M_i = m_i, \Omega_i = o_i) P(\Omega_i = o_i|A = a) \right. \\ & \left. \times \left\{ P(M_i = m_i|A = a + 1) - P(M_i = m_i|A = a) \right\} \right], \end{aligned} \quad (10)$$

where  $\mathbb{M}_i$  is the support of  $M_i$ ,  $o_i = [m_1, \dots, m_{i-1}, m_{i+1}, \dots, m_p]$ ,  $P(M = m|A = a)$  is the probability of  $M = m$  when setting  $A = a$ .

**Theorem 1.4** *Under assumptions (A1-A3) and Model 1, we have*

$$DM_i = \xi_i,$$

**Remark 1.5** *The proof can be found in Section 2.5. Based on Definition 3.2 and Equation 10, both the proposed  $DM$  and the effect defined in Vansteelandt and Daniel [5] contain the information of the causal effect of  $A$  on the mediator  $M_i$ , i.e.,  $P(M_i = m_i|A = a + 1) - P(M_i = m_i|A = a)$ .*

## 2 Technical Proofs

### 2.1 Proof of Theorem 3.1

**Proof 2.1** *In this proof, we will give the explicit expressions of causal effects defined under the LSEM. First, Equation 3 is equivalent to*

$$\begin{cases} A \equiv \epsilon_A, \\ M = \alpha A + B_M^\top M + \epsilon_M, \\ Y = \gamma A + \beta^\top M + \epsilon_Y. \end{cases} \quad (11)$$

Based on  $M = \alpha A + B_M^\top M + \epsilon_M$ , by moving  $B_M^\top M$  to the left-hand side, we have

$$(I_p - B_M^\top)M = \alpha A + \epsilon_M.$$

Suppose the mediators are sorted in the topological order (a series of elementary transformation of the matrix), then the matrix  $B_M^\top$  is strictly upper triangular with the diagonal element as 0. Thus, we have  $I_p - B_M^\top$  is invertible, then  $I_p - B_M^\top$  under its original order should be also invertible (any invertible matrix after elementary transformation is still invertible).

Therefore, we can rewrite  $M$  as a purely function of  $A$  plus the error term as follows.

$$M = (I_p - B_M^\top)^{-1} \alpha A + (I_p - B_M^\top)^{-1} \epsilon_M. \quad (12)$$

Then we replace mediators in Equation 11 with Equation 12 and obtain

$$\begin{cases} A & \equiv \epsilon_A, \\ M & = (I_p - B_M^\top)^{-1} \alpha A + (I_p - B_M^\top)^{-1} \epsilon_M, \\ Y & = \gamma A + \beta^\top M + \epsilon_Y \\ & = \gamma A + \{\beta^\top (I_p - B_M^\top)^{-1} \alpha\} A + \{\beta^\top (I_p - B_M^\top)^{-1} \epsilon_M + \epsilon_Y\}. \end{cases} \quad (13)$$

Next, we show how to get the explicit expressions of  $E\{Y|do(A=a)\}$  under the LSEM. Following the results in Rosenbaum and Rubin [4], under the assumption (A2), we have  $P\{M|do(A=a)\} = P(M|A=a)$ , and thus,

$$E\{M|do(A=a)\} = E(M|A=a).$$

Similarly, we can get  $E\{Y|do(A=a)\} = E(Y|A=a)$  under the assumption (A1), and  $E\{Y|do(A=a, M=m)\} = E(Y|A=a, M=m)$  under the assumption (A3).

Based on above results and Equation 13, we have

$$\begin{aligned} E\{Y|do(A=a)\} &= E\{Y|A=a\} \\ &= E\{\gamma A + \beta^\top M + \epsilon_Y|A=a\} \\ &= \gamma a + \beta^\top E\{M|A=a\} \\ &= \gamma a + \beta^\top E\{(I_p - B_M^\top)^{-1} \alpha A + (I_p - B_M^\top)^{-1} \epsilon_M|A=a\} \\ &= \gamma a + \beta^\top (I_p - B_M^\top)^{-1} \alpha a, \end{aligned} \quad (14)$$

where the first '=' is held under the assumption (A1), the second and forth '=' are given by Equation 13 that  $Y = \gamma A + \beta^\top M + \epsilon_Y$  and  $M = (I_p - B_M^\top)^{-1} \alpha A + (I_p - B_M^\top)^{-1} \epsilon_M$ .

Following the same calculation procedure of  $E\{Y|do(A = a)\}$ , we next give the natural direct effect under assumptions (A1-A3) and Model 1 as

$$\begin{aligned} DE &= E\{Y|do(A = a + 1, M = m^{(a)})\} - E\{Y|do(A = a)\} \\ &= \{\gamma(a + 1) + \beta^\top m^{(a)}\} - \{\gamma a + \beta^\top m^{(a)}\} \\ &= \gamma, \end{aligned}$$

where the first '=' is given by the definition of the DE.

Similarly, the natural indirect effect is

$$\begin{aligned} IE &= E\{Y|do(A = a, M = m^{(a+1)})\} - E\{Y|do(A = a)\} \\ &= \{\gamma a + \beta^\top m^{(a+1)}\} - \{\gamma a + \beta^\top m^{(a)}\} \\ &= \beta^\top (I_p - B_M^\top)^{-1} \alpha (a + 1) - \beta^\top (I_p - B_M^\top)^{-1} \alpha a \\ &= \beta^\top (I_p - B_M^\top)^{-1} \alpha. \end{aligned}$$

Thus, the total effect of  $A$  on  $Y$  is

$$TE = E\{Y|do(A = a + 1)\} - E\{Y|do(A = a)\} = DE + IE = \gamma + \beta^\top (I_p - B_M^\top)^{-1} \alpha.$$

Finally, we give the expression for the natural direct effect of  $M_i$  on  $Y$  under the LSEM. Based on the assumption (A2) and Equation 12, we have

$$\begin{aligned} &E\{M_i|do(A = a + 1)\} - E\{M_i|do(A = a)\} \\ &= E\{M_i|A = a + 1\} - E\{M_i|A = a\} \\ &= \{(I_p - B_M^\top)^{-1} \alpha\}_i (a + 1) - \{(I_p - B_M^\top)^{-1} \alpha\}_i a \\ &= \{(I_p - B_M^\top)^{-1} \alpha\}_i, \end{aligned} \tag{15}$$

where  $\{(I_p - B_M^\top)^{-1} \alpha\}_i$  is the  $i$ -th element of the vector  $(I_p - B_M^\top)^{-1} \alpha$ .

Then, based on  $Y = \gamma A + \beta^\top M + \epsilon_Y$  and the assumption (A3), we have,

$$\begin{aligned} &E\{Y|do(A = a, M_i = m_i^{(a)} + 1, \Omega_i = o_i^{(a)})\} - E\{Y|do(A = a)\} \\ &= E\{Y|A = a, M_i = m_i^{(a)} + 1, \Omega_i = o_i^{(a)}\} - E\{Y|A = a\} \\ &= \gamma a + \beta^\top \begin{bmatrix} m_1^{(a)} \\ \vdots \\ m_i^{(a)} + 1 \\ \vdots \\ m_p^{(a)} \end{bmatrix} - \gamma a - \beta^\top \begin{bmatrix} m_1^{(a)} \\ \vdots \\ m_i^{(a)} \\ \vdots \\ m_p^{(a)} \end{bmatrix} = \beta^\top \mathbf{1}_i = \beta_i, \end{aligned}$$

where  $\mathbf{1}_i$  is a  $p \times 1$  vector with the  $i$ -th element as 1 while others equal to 0, and  $\beta_i$  is the  $i$ -th element of the vector  $\beta$ .

Thus, we have

$$\begin{aligned}
DM_i &= \left[ E\{M_i|do(A = a + 1)\} - E\{M_i|do(A = a)\} \right] \\
&\quad \times \left[ E\{Y|do(A = a, M_i = m_i^{(a)} + 1, \Omega_i = o_i^{(a)})\} - E\{Y|do(A = a)\} \right], \\
&= \{(I_p - B_M^\top)^{-1} \boldsymbol{\alpha}\}_i \times \beta_i \\
&= \beta_i \{(I_p - B_M^\top)^{-1} \boldsymbol{\alpha}\}_i. \quad \square
\end{aligned}$$

## 2.2 Proof of Theorem 1.1

**Proof 2.2** 1. If there is no directed path  $\pi^* \in \{\pi_{AY}(\mathcal{G})\}$  such that the length of  $\pi^*$  is larger than 2, i.e., the length of  $\pi^* \in \{\pi_{AY}(\mathcal{G})\}$  is either 1 or 2. Here, the path with length 1 corresponds to  $A \rightarrow Y$ , and paths with length 2 are  $A \rightarrow M_i \rightarrow Y$  with possibly  $i = 1, \dots, p$ . Thus, there is no interaction among mediators.

By the definition of the LSEM, we have  $B_M = \mathbf{0}_{p \times p}$ , where  $\mathbf{0}_{p \times p}$  is a  $p \times p$  zero matrix. Following the path method (the causal effect of  $X_i$  on  $X_j$  along a directed path from  $X_i \rightarrow X_j$  in  $\mathcal{G}$  can be calculated by multiplying all edge weights along the path) illustrated in Wright [6] and Nandy et al. [3], we could obtain  $\sum \eta_i = \sum_i \beta_i \boldsymbol{\alpha}_i = IE$ . (See a toy example provided in section ?? to illustrate how to use the path method to manually compute the causal effects.)

2. If there exists at least one directed path  $\pi^* \in \{\pi_{AY}(\mathcal{G})\}$  such that the length of  $\pi^*$  is larger than 2, and the element in  $B$  is nonnegative, we have  $B_M \neq \mathbf{0}_{p \times p}$ . Without loss of generality, suppose there exists  $M_i \in M$  with a set of directed path that starts with  $A$ , contains  $M_i$ , then goes through other mediators, and ends with  $Y$ , denoted each path in such set as  $\pi_{i,j} = \{A \rightarrow \dots \rightarrow M_i \dots \rightarrow \dots \rightarrow Y\}$  for  $j = 1, \dots, n_i$ , where  $n_i$  is the size of such path set for  $M_i$ , and the weights of edges in  $\pi_{i,j}$  is positive. Note the set  $\{\pi_{i,j}\}$  excludes the paths end with  $M_i \rightarrow Y$ .

Let  $e_{\pi_{i,j}}$  denote the causal effect of  $A$  on  $Y$  through directed path  $\pi_{i,j}$ . Based on the path method in Wright [6] and Nandy et al. [3] with the definition of  $IM_i$ , we have its theoretical form as

$$IM_i = \sum_{j=1}^{n_i} e_{\pi_{i,j}}. \quad (16)$$

By Equation 15 and the definition of  $\eta_i$ , we have its first multiplier as

$$E\{M_i|do(A = a + 1)\} - E\{M_i|do(A = a)\} = \{(I_p - B_M^\top)^{-1} \boldsymbol{\alpha}\}_i,$$

which is also the first multiplier in both  $DM_i$  and  $IM_i$ .

And the second multiplier of  $\eta_i$  can be expressed as

$$\begin{aligned}
&E\{Y|do(M_i = m_i + 1)\} - E\{Y|do(M_i = m_i)\} \\
&= E\{Y|do(M_i = m_i^{(a)} + 1)\} - E\{Y|do(M_i = m_i^{(a)})\} \\
&= E\{Y|do(A = a, M_i = m_i^{(a)} + 1)\} - E\{Y|do(A = a, M_i = m_i^{(a)})\}, \\
&= E\{Y|do(A = a, M_i = m_i^{(a)} + 1)\} - E\{Y|do(A = a)\},
\end{aligned}$$

where  $m_i^{(a)}$  is the value of  $M_i$  when setting  $do(A = a)$ . Here, the first '=' is valid since  $m_i$  can be arbitrary number, and the second and third '=' are based on the equivalent interventions.

Based on the technique of plus and minus the same term, we decompose the second multiplier of  $\eta_i$  into two parts as follows

$$\begin{aligned}
& E\{Y|do(A = a, M_i = m_i^{(a)} + 1)\} - E\{Y|do(A = a)\} \\
&= \underbrace{\left[ E\{Y|do(A = a, M_i = m_i^{(a)} + 1, \Omega_i = o_i^{(a)})\} - E\{Y|do(A = a)\} \right]}_{\text{the second multiplier of } DM_i} \\
&+ \underbrace{\left[ E\{Y|do(A = a, M_i = m_i^{(a)} + 1)\} - E\{Y|do(A = a, M_i = m_i^{(a)} + 1, \Omega_i = o_i^{(a)})\} \right]}_{\text{the second multiplier of } IM_i}
\end{aligned} \tag{17}$$

where  $\Omega_i = M \setminus M_i$  is the sets of mediators except  $M_i$ , and  $o_i^{(a)}$  is the value of  $\Omega_i$  when setting  $do(A = a)$ . Here, the first term in the above equation corresponds to the second multiplier of  $DM_i$ , while the second term is the second multiplier of  $IM_i$ .

Thus, the summation of  $\eta_i$  is

$$\begin{aligned}
\sum \eta_i &= \sum_i \left\{ \left[ E\{M_i|do(A = a + 1)\} - E\{M_i|do(A = a)\} \right] \right. \\
&\quad \left. \times \left[ E\{Y|do(M_i = m_i + 1)\} - E\{Y|do(M_i = m_i)\} \right] \right\} \\
&= \sum_i \{DM_i + IM_i\} = \sum_i DM_i + \sum_i IM_i = IE + \sum_i \sum_{j=1}^{n_i} e_{\pi_{i,j}},
\end{aligned}$$

where the first '=' is from Definition 1, the second '=' is given by Equation 17 and Definition 3.2 and 3.3, and the last '=' comes from Theorem 3.2 and the theoretical form of  $IM$  in Equation 16.

Here, we have  $e_{\pi_{i,j}} > 0$  since the weights of edges in  $\pi_{i,j}$  is positive based on the path method in Wright [6] and Nandy et al. [3]. Then,  $\sum_i \sum_{j=1}^{n_i} e_{\pi_{i,j}}$  is also strictly larger than 0. Therefore, we have

$$\sum \eta_i > IE. \quad \square$$

## 2.3 Proof of Theorem 1.2

**Proof 2.3** 1. If there doesn't exist edge  $M_i \rightarrow Y$  in  $\mathcal{G}$ , then by definition we have  $\beta_i = 0$ . Thus,  $LE_i = \beta_i \zeta_i = 0$ .

2. If there exists edge  $M_i \rightarrow Y$  in  $\mathcal{G}$ . Suppose there is a directed path set with size  $m_i$  associated to the edge  $M_i \rightarrow Y$ , where each directed path  $\tilde{\pi}_{i,j}$  starts with node  $A$  and ends with  $M_i \rightarrow Y$ , denoted as  $\tilde{\pi}_{i,j} = \{A \rightarrow \dots \rightarrow \dots \rightarrow M_i \rightarrow Y\}$  for  $j = 1, \dots, m_i$ .

Let  $e_{\tilde{\pi}_{i,j}}$  denote the causal effect of  $A$  on  $Y$  through directed path  $\tilde{\pi}_{i,j}$ ,  $e_{\tilde{\pi}_{i,j}}^{(A, M_i)}$  be the causal effect of  $A$  on  $M_i$  through directed path  $\tilde{\pi}_{i,j}$ , and  $e^{(M_i, Y)}$  is the causal effect of  $M_i$  on  $Y$  through

directed edge  $M_i \rightarrow Y$ . Following the path method in Wright [6] and Nandy et al. [3], we have  $e_{\tilde{\pi}_{i,j}} = e_{\tilde{\pi}_{i,j}}^{(A,M_i)} e^{(M_i,Y)}$ .

Then the  $i$ -th last edge-specific effect is equal to the summation of the effect through each path  $\tilde{\pi}_{i,j}$ , i.e.,

$$LE_i = \sum_{j=1}^{n_i} e_{\tilde{\pi}_{i,j}} = \sum_{j=1}^{n_i} e_{\tilde{\pi}_{i,j}}^{(A,M_i)} e^{(M_i,Y)} = e^{(M_i,Y)} \sum_{j=1}^{n_i} e_{\tilde{\pi}_{i,j}}^{(A,M_i)}.$$

Here, by the similar argument based on the path method, we have  $e^{(M_i,Y)} = \beta_i$  and  $\sum_{j=1}^{n_i} e_{\tilde{\pi}_{i,j}}^{(A,M_i)}$  as the total causal effect of  $A$  on  $M_i$ .

Recall that  $\zeta_i \equiv \{(I_p - B_M^\top)^{-1} \alpha\}_i$  is the causal effect of  $A$  on  $M_i$ . Therefore, the  $i$ -th LE is the product of the causal effect of  $A$  on  $M_i$  and the causal effect of  $M_i$  on  $Y$ , i.e.,

$$LE_i = \beta_i \zeta_i. \quad \square$$

## 2.4 Proof of Theorem 1.3

**Proof 2.4** Given a general DAG  $\mathcal{G}$  with nodes  $\{A, M, Y\}$ , let the union of all directed paths that contain the  $i$ -th last edge as  $\tau_i = \{\pi : A \rightarrow \dots \rightarrow M_i \rightarrow Y\}, i = 1, \dots, p$ . Here, we have  $\tau_i = \{\tilde{\pi}_{i,j}\}_{1 \leq j \leq m_j}$  established in Section 2.3. It is clear that the union set of  $\tau_i$  in  $\mathcal{G}$  is equal to the set of all directed paths start with  $A$  and end with node  $Y$  (except  $A \rightarrow Y$ ) in  $\mathcal{G}$  as

$$\bigcup_i \tau_i = \{\pi_{AY}(\mathcal{G})\} \setminus \{A \rightarrow Y\}.$$

Also, based on the uniqueness of each last edge,  $\tau_i$  is pairwise disjoint, i.e.

$$\tau_i \cap \tau_j = \emptyset, \quad \forall i \neq j.$$

Since the IE is defined as the total causal effect of  $A$  on  $Y$  that goes through mediators, we have the IE equal to the causal effect that goes through the set  $\{\pi_{AY}(\mathcal{G})\} \setminus \{A \rightarrow Y\}$ , i.e., the IE equal to the causal effect that goes through set  $\bigcup_i \tau_i$ . Based on the mutual disjoint property of  $\tau_i$ , we have the IE is exactly the summation of the causal effect through  $\tau_i$ . Lastly, from the definition of  $LE_i$ , we have

$$IE = \sum_{i=1}^p LE_i. \quad \square$$

## 2.5 Proof of Theorem 1.4

**Proof 2.5** The proof of the consistency of our defined DM to the interventional effect  $\xi_i$  can be completed based on Equation 3 under assumptions (A1-A3) and Model 1.



Recall the definition in Equation 10, we have

$$\begin{aligned}
\xi_i &= \sum_{m_1 \in \mathbb{M}_1} \cdots \sum_{m_p \in \mathbb{M}_p} \left[ E(Y|A = a, M_i = m_i, \Omega_i = o_i) P(\Omega_i = o_i|A = a) \right. \\
&\quad \left. \times \left\{ P(M_i = m_i|A = a + 1) - P(M_i = m_i|A = a) \right\} \right] \\
&= \sum_{m_1 \in \mathbb{M}_1} \cdots \sum_{m_p \in \mathbb{M}_p} \left\{ E(Y|A = a, M_i = m_i, \Omega_i = o_i) P(\Omega_i = o_i|A = a) P(M_i = m_i|A = a + 1) \right. \\
&\quad \left. - E(Y|A = a, M_i = m_i, \Omega_i = o_i) P(\Omega_i = o_i|A = a) P(M_i = m_i|A = a) \right\}.
\end{aligned}$$

Given  $A = a$ , the value of  $M_i$  is  $m_i^{(a)}$  and  $\Omega_i$  takes  $o_i^{(a)}$ ; while when setting  $A = a + 1$ , the value of  $M_i$  is  $m_i^{(i)}$ . Therefore, we have  $P(M_i = m_i|A = a) = 1$  if  $m_i = m_i^{(a)}$  otherwise is 0, and  $P(\Omega_i = o_i|A = a) = 1$  if  $o_i = o_i^{(a)}$  otherwise is 0.

Under assumptions (A1-A3), we have

$$\xi_i = E(Y|A = a, M_i = m_i^{(a+1)}, \Omega_i = o_i^{(a)}) - E(Y|A = a, M_i = m_i^{(a)}, \Omega_i = o_i^{(a)}).$$

Then, based on the LSEM that  $Y = \gamma A + \beta^\top M + \epsilon_Y$ , we can further obtain that

$$\xi_i = \gamma a + \beta^\top \begin{bmatrix} m_1^{(a)} \\ \vdots \\ m_i^{(a+1)} \\ \vdots \\ m_p^{(a)} \end{bmatrix} - \gamma a - \beta^\top \begin{bmatrix} m_1^{(a)} \\ \vdots \\ m_i^{(a)} \\ \vdots \\ m_p^{(a)} \end{bmatrix} = \beta_i \{m_i^{(a+1)} - m_i^{(a)}\}.$$

From Equation 12, we have

$$\xi_i = \beta_i \left[ \{(I_p - B_M^\top)^{-1} \alpha\}_i (a + 1) - \{(I_p - B_M^\top)^{-1} \alpha\}_i a \right] = \beta_i \{(I_p - B_M^\top)^{-1} \alpha\}_i.$$

Thus, under assumptions (A1-A3) and Model 1, we have

$$DM_i = \xi_i. \quad \square$$

## References

- [1] Avin, C., Shpitser, I. and Pearl, J. [2005], ‘Identifiability of path-specific effects’.
- [2] Chakraborty, A., Nandy, P. and Li, H. [2018], ‘Inference for individual mediation effects and interventional effects in sparse high-dimensional causal graphical models’, *arXiv preprint arXiv:1809.10652*.

- [3] Nandy, P., Maathuis, M. H., Richardson, T. S. et al. [2017], ‘Estimating the effect of joint interventions from observational data in sparse high-dimensional settings’, *The Annals of Statistics* **45**(2), 647–674.
- [4] Rosenbaum, P. R. and Rubin, D. B. [1983], ‘The central role of the propensity score in observational studies for causal effects’, *Biometrika* **70**(1), 41–55.
- [5] Vansteelandt, S. and Daniel, R. M. [2017], ‘Interventional effects for mediation analysis with multiple mediators’, *Epidemiology (Cambridge, Mass.)* **28**(2), 258.
- [6] Wright, S. [1921], ‘Correlation and causation’, *Journal of agricultural research* **20**(7), 557–585.