

# The Principal Direction of Mediation Supplementary Materials

Oliver Chén, Elizabeth Ogburn, Ciprian Crainiceanu,  
Brian Caffo, Martin Lindquist

## 1 Proofs

### Lemma 1 (Consistency Theorem):

Suppose that  $Q(\boldsymbol{\theta}; \mathbf{Z}_n)$  is continuous in  $\boldsymbol{\theta}$  and there exists a function  $Q_0(\boldsymbol{\theta})$  such that:

- (1)  $Q_0(\boldsymbol{\theta})$  is uniquely maximized at  $\boldsymbol{\theta}_0$ ;
- (2)  $\Theta$  is compact;
- (3)  $Q_0(\boldsymbol{\theta})$  is continuous in  $\boldsymbol{\theta}$ ;
- (4)  $Q(\boldsymbol{\theta}; \mathbf{Z}_n)$  converges uniformly in probability to  $Q_0(\boldsymbol{\theta})$ .

Then  $\hat{\boldsymbol{\theta}}(\mathbf{Z}_n)$  defined as the value of  $\boldsymbol{\theta} \in \Theta$  which for each  $\mathbf{Z}_n = \mathbf{z}_n$  maximizes the objective function  $Q(\boldsymbol{\theta}; \mathbf{Z}_n)$  satisfies  $\hat{\boldsymbol{\theta}}(\mathbf{Z}_n) \xrightarrow{p} \boldsymbol{\theta}_0$ .

### Proof of Lemma 1:

See Theorem 2.1 in [\(Newey and McFadden, 1994\)](#). ■

### Lemma 2

Consider a compact space  $\Theta$ . Consider:

$$L(z, \mathbf{w}, \boldsymbol{\theta}, \lambda) = f^{\text{obj}}(z, \mathbf{w}, \boldsymbol{\theta}) + f^{\text{pen}}(z, \mathbf{w}, \lambda),$$

where  $f^{\text{obj}}(z, \mathbf{w}, \boldsymbol{\theta})$  is an objective function and  $f^{\text{pen}}(z, \mathbf{w}, \lambda) = \frac{\lambda\{f^{\text{cons}}(\mathbf{w}) - c\}}{n}$  is a penalization function.

If:

- (1) both the objective function and the penalization function can be profiled by  $\theta$ , defined as  $f^{\text{obj}}(z, \theta)$  and  $f^{\text{pen}}(\theta) := \frac{\lambda(\theta)\{f^{\text{cons}}(\theta) - c\}}{n}$ ;
- (2) the objective function is a log likelihood function;
- (3) both  $f^{\text{obj}}(z, \theta)$  and  $f^{\text{pen}}(\theta)$  are continuous in  $\theta$ ;
- (4) if there exists a function  $d_0(z)$  such that  $|L(z, \theta)| := |f^{\text{obj}}(z, \theta) + f^{\text{pen}}(\theta)| \leq d_0(z)$  for all  $\theta \in \Theta$  and  $z \in \mathcal{Z}$ , and  $\mathbb{E}_{\theta_0}[d_0(x)] < \infty$ , then:

- i.  $q_0(\theta) := \mathbb{E}_{\theta_0}[L(z, \theta)]$  is continuous in  $\theta$ ;
- ii.  $\sup_{\theta \in \Theta} |q(\theta; \mathbf{Z}_n) - q_0(\theta)| \xrightarrow{p} 0$ , where  $q(\theta; \mathbf{Z}_n) := \frac{1}{n}L(\mathbf{Z}_i, \theta)$ .

Note: the above Lemma can be stated in a more general case where there are multiple sets of parameters and several constraint functions.

### **Proof of Lemma 2:**

Consider the regularity conditions in the Appendix.

#### **Part i:**

$\forall \theta \in \Theta$ , choose a sequence  $\theta_k \in \Theta$ , such that  $\theta_k \rightarrow \theta$ . By (N-3), we have  $L(x; \theta_k) \rightarrow L(x; \theta)$ . By (N-4), then  $q_0(\theta_k) := \mathbb{E}_{\theta_0}(L(Z, \theta_k)) \rightarrow \mathbb{E}_{\theta_0}(L(Z, \theta)) = q_0(\theta)$ , by the dominated convergence theorem (DCT henceforth). Hence,  $q_0(\theta)$  is continuous in  $\theta$ .

#### **Part ii:**

We need to show that  $\forall \epsilon, \xi, \exists N(\epsilon, \xi)$  such that  $\forall n > N(\epsilon, \xi)$ ,

$$P[\sup_{\theta \in \Theta} |q(\theta; \mathbf{Z}_n) - q_0(\theta)| < \xi].$$

Since both  $f^{\text{obj}}(z, \theta)$  and  $f^{\text{prof}}(\theta)$  are continuous in  $\theta$ , then  $L(z, \theta)$  is uniformly continuous.

Hence,

$$\Delta(z, \delta) = \sup_{\{(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2) : \|\boldsymbol{\theta}_1 - \boldsymbol{\theta}_2\| < \delta\}} |L(z, \boldsymbol{\theta}_1) - L(z, \boldsymbol{\theta}_2)| \rightarrow 0$$

as  $\delta \rightarrow 0$ .

By (N-4),  $\Delta(z, \delta) \leq 2d_0(z) \forall \delta$ .

By DCT,  $\mathbb{E}_{\boldsymbol{\theta}_0}[\Delta(Z, \delta)] \rightarrow 0$  as  $\delta \rightarrow 0$ .

Define  $B(\boldsymbol{\theta}_j, \delta) = \{\tilde{\boldsymbol{\theta}} : \|\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_j\| < \delta\}$ . Since  $\Theta$  is compact, for every fixed  $\delta$ ,  $\exists$  a subcover  $\{B(\boldsymbol{\theta}_j, \delta), j = 1, \dots, J\}$  such that  $\bigcup_{j=1}^{J < \infty} B(\boldsymbol{\theta}_j, \delta) \supset \Theta$ .

By the triangle inequality, we have:

$$|q(\boldsymbol{\theta}; \mathbf{Z}_n) - q_0(\boldsymbol{\theta})| \leq |q(\boldsymbol{\theta}; \mathbf{Z}_n) - q(\boldsymbol{\theta}_j; \mathbf{Z}_n)| \quad (1)$$

$$+ |q(\boldsymbol{\theta}_j; \mathbf{Z}_n) - q_0(\boldsymbol{\theta}_j)| \quad (2)$$

$$+ |q_0(\boldsymbol{\theta}_j) - q_0(\boldsymbol{\theta})|. \quad (3)$$

Choose  $\boldsymbol{\theta}_j$  such that  $\boldsymbol{\theta} \in B(\boldsymbol{\theta}_j; \delta)$ . Since  $\|\boldsymbol{\theta} - \boldsymbol{\theta}_j\| < \delta$ , then:

$$(1) = \left| \frac{1}{n} \sum_{i=1}^n \{L(Z_i, \boldsymbol{\theta}) - L(Z_i, \boldsymbol{\theta}_j)\} \right| \leq \frac{1}{n} \sum_{i=1}^n |L(Z_i, \boldsymbol{\theta}) - L(Z_i, \boldsymbol{\theta}_j)| \leq \frac{1}{n} \sum_{i=1}^n \Delta(Z_i, \delta).$$

Next,

$$(2) < \max_{j \in \{1, \dots, J\}} |q(\boldsymbol{\theta}_j; \mathbf{Z}_n) - q_0(\boldsymbol{\theta}_j)|$$

Moreover, choosing  $\delta$  to be small,

$$(3) \leq \sup_{\{(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2) : \|\boldsymbol{\theta}_1 - \boldsymbol{\theta}_2\| < \delta\}} |q_0(\boldsymbol{\theta}_1) - q_0(\boldsymbol{\theta}_2)| \leq \epsilon^*(\delta)$$

where  $\epsilon(\delta) \rightarrow 0$  as  $\delta \rightarrow 0$ .

Combine (1) - (3), we have:

$$\sup_{\boldsymbol{\theta} \in \Theta} |q(\boldsymbol{\theta}; \mathbf{Z}_n) - q_0(\boldsymbol{\theta})| \leq \frac{1}{n} \sum_{i=1}^n \Delta(Z_i, \delta) + \max_{j \in \{1, \dots, J\}} |q(\boldsymbol{\theta}_j; \mathbf{Z}_n) - q_0(\boldsymbol{\theta}_j)| + \epsilon^*(\delta).$$

Choose  $\delta$  such that  $\epsilon^*(\delta) < \frac{\epsilon}{3}$ . Call this  $\delta_1$ . So for any  $\delta < \delta_1$ , we have:

$$\begin{aligned}
& P_{\theta_0}[\sup_{\theta \in \Theta} |q(\theta; \mathbf{Z}_n) - q_0(\theta)| > \epsilon] \\
& \leq P_{\theta_0}[\frac{1}{n} \sum_{i=1}^n \Delta(Z_i, \delta) + \max_{j \in \{1, \dots, J\}} |q(\theta_j; \mathbf{Z}_n) - q_0(\theta_j)| > \frac{2\epsilon}{3}] \\
& \leq P_{\theta_0}[\frac{1}{n} \sum_{i=1}^n \Delta(Z_i, \delta) > \frac{\epsilon}{3}] + \tag{4}
\end{aligned}$$

$$P_{\theta_0}[\max_{j \in \{1, \dots, J\}} |q(\theta_j; \mathbf{Z}_n) - q_0(\theta_j)| > \frac{\epsilon}{3}] \tag{5}$$

Note (4) =  $P_{\theta_0}[\frac{1}{n} \sum_{i=1}^n \{\Delta(Z_i, \delta) - \mathbb{E}_{\theta_0}[\Delta(Z; \delta)]\} + \mathbb{E}_{\theta_0}[\Delta(Z; \delta)] > \frac{\epsilon}{3}]$ , where  $\mathbb{E}_{\theta_0}[\Delta(Z; \delta)] \rightarrow 0$  as  $\delta \rightarrow 0$ .

Choose  $\delta$  small enough such that  $\mathbb{E}_{\theta_0}[\Delta(Z; \delta)] < \frac{\epsilon}{6}$ . Call this  $\delta_2$ .

Take  $\delta < \min(\delta_1, \delta_2)$ . Then:

$$P_{\theta_0}[\frac{1}{n} \sum_{i=1}^n \{\Delta(Z_i, \delta) - \mathbb{E}_{\theta_0}[\Delta(Z; \delta)]\} > \frac{\epsilon}{6}] := (4)'$$

By the Weak Law of Large Numbers (henceforth WLLN),  $\exists N_1(\epsilon, \xi)$  such that  $\forall n > N_1(\epsilon, \xi)$ ,  $(4) < (4)' < \frac{\xi}{2}$ .

Consider the finite subcover  $\{B(\theta_j, \delta), j = 1, \dots, J\}$  for  $\delta$  above considered. Note:

$$\begin{aligned}
(5) &= P_{\theta_0}[\bigcup_{j=1}^J \{|q(\theta_j; \mathbf{Z}_n) - q_0(\theta_j)| > \frac{\epsilon}{3}\}] \\
&\leq \sum_{j=1}^J P_{\theta_0}[|q(\theta_j; \mathbf{Z}_n) - q_0(\theta_j)| > \frac{\epsilon}{3}].
\end{aligned}$$

By the WLLN,  $\forall \theta_j$  and  $\forall \epsilon, \xi > 0$ ,  $\exists N_{2j}(\epsilon, \xi)$  such that  $\forall n > N_{2j}(\epsilon, \xi)$ :

$$P_{\theta_0}[|q(\theta_j; \mathbf{Z}_n) - q_0(\theta_j)| > \frac{\epsilon}{3}] \leq \frac{\xi}{2J}.$$

Let  $N_2(\epsilon, \xi) = \max_{j \in \{1, \dots, J\}} \{N_{2j}\}$ . Then,  $\forall n > N_2(\epsilon, \xi)$ , we have:

$$\sum_{j=1}^J P_{\theta_0}[|q(\theta_j; \mathbf{Z}_n) - q_0(\theta_j)| > \frac{\epsilon}{3}] < \frac{\xi}{2}.$$

Hence, (5)  $< \frac{\xi}{2}$ .

The result for (4) and (5) show that  $\exists$  an  $N(\epsilon, \xi) = \max(N_1(\epsilon, \xi), N_2(\epsilon, \xi))$  such that  $\forall n > N(\epsilon, \xi)$ ,

$$P_{\theta_0}[\sup_{\theta \in H} | q(\theta; \mathbf{Z}_n) - q_0(\theta) |] < \xi. \blacksquare$$

**Proof of Theorem 1:**

Define  $Q(\theta; \mathbf{Z}_n) := \frac{1}{n} L(\theta; \mathbf{Z}_n)$  and  $Q_0(\theta) := \mathbb{E}_{\theta_0}(L(\theta; \mathbf{Z}_n))$ .

Recall  $L(\theta; z) := f^{\text{obj}}(\theta; z) + \frac{\lambda(\theta)(f^{\text{cons}} - c)}{n}$ .

**(i. Consistency)**

We want to show  $\hat{\theta}(\mathbf{Z}_n) := \operatorname{argmax}_{\theta} \{L(\mathbf{Z}_n; \theta)\}$  converges to  $\theta_0$  in probability, i.e.

$$\hat{\theta} \xrightarrow{p} \theta_0. \tag{6}$$

By Lemma 1, it suffices to show:

- (a)  $Q(\theta; \mathbf{Z}_n)$  is continuous in  $\theta$ ;
- (b)  $Q_0(\theta)$  is continuous in  $\theta$ ;
- (c)  $\sup_{\theta \in H} | Q(\theta; \mathbf{Z}_n) - Q_0(\theta) | \xrightarrow{p} 0$ .

(a) is implied by (N-3);

Due to (N-3) and (N-4), by Lemma 2, we have:

- (d)  $q_0(\theta) := \mathbb{E}_{\theta_0}(L(z, \theta))$  is continuous in  $\theta$ ; and
- (e)  $\sup_{\theta \in H} \| \hat{q}(\theta; \mathbf{Z}_n) - q_0(\theta) \| \xrightarrow{p} 0$ , where  $\hat{q}(\theta; \mathbf{Z}_n) = \frac{1}{n} \sum_{i=1}^n L(Z_i, \theta)$
- (d) implies (b);
- (e) implies (c).  $\blacksquare$

**(ii. Asymptotic Normality)**

The estimator  $\hat{\theta}(\mathbf{Z}_n)$  (henceforth  $\hat{\theta}$ ) satisfies:

$$\begin{aligned} \dot{L}(\mathbf{Z}_n, \hat{\boldsymbol{\theta}}) &= \frac{\partial L(\mathbf{Z}_n, \boldsymbol{\theta})}{\partial \boldsymbol{\theta}} /_{\boldsymbol{\theta}=\hat{\boldsymbol{\theta}}} \\ &= \sum_{i=1}^n \left\{ \frac{\partial f^{\text{obj}}(Z_i; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}} + \frac{\nabla^{\boldsymbol{\theta}} \lambda(\boldsymbol{\theta})(f^{\text{cons}}(\boldsymbol{\theta}) - c) + \lambda(\boldsymbol{\theta}) \nabla^{\boldsymbol{\theta}} f^{\text{cons}}(\boldsymbol{\theta})}{n} \right\} = 0. \end{aligned}$$

Recall that we have defined:

$$\ell(z; \boldsymbol{\theta}) = \frac{\partial f^{\text{obj}}(z; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}} + \frac{\nabla^{\boldsymbol{\theta}} \lambda(\boldsymbol{\theta})(f^{\text{cons}}(\boldsymbol{\theta}) - c) + \lambda(\boldsymbol{\theta}) \nabla^{\boldsymbol{\theta}} f^{\text{cons}}(\boldsymbol{\theta})}{n}.$$

Define:

$$\hat{q}(\mathbf{Z}_n; \boldsymbol{\theta}) = \frac{1}{n} \sum_{i=1}^n L(Z_i; \boldsymbol{\theta}),$$

$$q_0(\boldsymbol{\theta}) = \mathbb{E}_{\boldsymbol{\theta}_0}(L(z; \boldsymbol{\theta})),$$

$$\hat{q}(\mathbf{Z}_n; \boldsymbol{\theta}) = \frac{1}{n} \sum_{i=1}^n \ell(Z_i; \boldsymbol{\theta}),$$

$$\dot{q}_0(\boldsymbol{\theta}) = \mathbb{E}_{\boldsymbol{\theta}_0}(\ell(z; \boldsymbol{\theta})).$$

Hence,  $\hat{\boldsymbol{\theta}}$  satisfies

$$\hat{q}(\mathbf{Z}_n; \hat{\boldsymbol{\theta}}) = 0. \quad (7)$$

By (N-2) and (N-6),

$$\dot{q}_0(\boldsymbol{\theta}) = O_p(n^{-1/2}). \quad (8)$$

Expanding  $\hat{q}(\mathbf{Z}_n; \hat{\boldsymbol{\theta}})$  around  $\boldsymbol{\theta}_0$ , we have:

$$0 = \hat{q}(\mathbf{Z}_n; \hat{\boldsymbol{\theta}}) = \hat{q}(\mathbf{Z}_n; \boldsymbol{\eta}_0) + D_n^*(\mathbf{Z}_n)(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0) \quad (9)$$

where  $D_n^*(\mathbf{Z}_n) = \hat{D}(\mathbf{Z}_n; \boldsymbol{\theta}^*) = \nabla^{\boldsymbol{\theta}} \hat{q}(\mathbf{Z}_n; \boldsymbol{\theta}) /_{\{\boldsymbol{\theta}=\boldsymbol{\theta}^*: \boldsymbol{\theta}^* \in \overline{(\boldsymbol{\theta}_0, \hat{\boldsymbol{\theta}})}\}}$ ,  $\nabla^{\boldsymbol{\theta}} \hat{q}(\mathbf{Z}_n; \boldsymbol{\theta}) = \frac{\partial \hat{q}(\mathbf{Z}_n; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}}$ , and  $\overline{(\boldsymbol{\theta}_0, \hat{\boldsymbol{\theta}})}$  denotes the open interval bounded by  $\boldsymbol{\theta}_0$  and  $\hat{\boldsymbol{\theta}}$ .

Assume (N-9) and (N-10), by Lemma 2, we have:

$$\sup \|\hat{D}(\mathbf{Z}_n; \boldsymbol{\theta}) - D_0(\boldsymbol{\theta})\| \xrightarrow{p} 0 \quad (10)$$

where  $D_0(\boldsymbol{\theta}) = \mathbb{E}_{\boldsymbol{\theta}_0}(\nabla^{\boldsymbol{\theta}} \ell(z; \boldsymbol{\theta}))$  and  $D_0(\boldsymbol{\theta}_0) = \mathbb{E}_{\boldsymbol{\theta}_0}(\nabla^{\boldsymbol{\theta}} \ell(z; \boldsymbol{\theta}_0))$ .

Since  $\boldsymbol{\theta}^* \in \mathcal{N}_r(\boldsymbol{\theta}_0)$  w.p.1.,

$$D_n^*(\mathbf{Z}_n) \xrightarrow{p} D_0(\boldsymbol{\theta}_0). \quad (11)$$

Hence, from (9), we have:

$$\sqrt{n}(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0) = -\{D_n^*(\mathbf{Z}_n)\}^{-1} \{\sqrt{n}\hat{q}(\mathbf{Z}_n; \boldsymbol{\theta}_0)\}. \quad (12)$$

By the C.L.T., the result from (8), and Slutsky's theorem, we have:

$$\sqrt{n}\hat{q}(\mathbf{Z}_n; \boldsymbol{\theta}_0) = \sqrt{n} \left( \frac{\sum_{i=1}^n \ell(Z_i; \boldsymbol{\theta}_0)}{n} - \dot{q}_0(\boldsymbol{\theta}) \right) + \sqrt{n}\dot{q}_0(\boldsymbol{\theta}) \rightarrow N(m(\boldsymbol{\theta}_0), V(\boldsymbol{\theta}_0)) \quad (13)$$

where  $V(\boldsymbol{\theta}_0) = \mathbb{E}_{\boldsymbol{\theta}_0}(\ell(z; \boldsymbol{\theta})\ell^{\top}(z; \boldsymbol{\theta})) - m(\boldsymbol{\theta}_0)[m(\boldsymbol{\theta}_0)]^{\top}$  and  $m(\boldsymbol{\theta}_0) = O_p(1)$ .

By (11) - (13), we have:

$$\sqrt{n}(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0) \rightarrow N(\mu(\boldsymbol{\theta}_0), \Sigma(\boldsymbol{\theta}_0)) \quad (14)$$

where  $\mu(\boldsymbol{\theta}_0) = D_0^{-1}(\boldsymbol{\theta}_0)m(\boldsymbol{\theta}_0)$ ,  $\Sigma(\boldsymbol{\theta}_0) = D_0^{-1}(\boldsymbol{\theta}_0)V(\boldsymbol{\theta}_0)[D_0^{-1}(\boldsymbol{\theta}_0)]^{\top}$ , and  $m(\boldsymbol{\theta}_0) = O_p(1)$ .

■

**Lemma 3** (Multivariate Delta Method When  $\overset{\circ}{J} > \overset{\circ}{J}$ , where  $g(\boldsymbol{\theta}) : \mathbb{R}^{\overset{\circ}{J}} \mapsto \mathbb{R}^{\overset{\circ}{J}}$ )

Many a time the Multivariate Delta Method is used without specifying dimension differences between the domain space and codomain space, or when it does, it is implied when the dimension of the domain space,  $\overset{\circ}{J}$ , is larger than or equal to it of the codomain space,  $\overset{\circ}{J}$ . Since our method and theory apply for both cases when  $\overset{\circ}{J} \geq \overset{\circ}{J}$  and  $\overset{\circ}{J} < \overset{\circ}{J}$ , for completeness purpose, we include a proof of the Multivariate Delta Method with a particular emphasis that the method is applicable for cases when  $\overset{\circ}{J} < \overset{\circ}{J}$ . The  $\overset{\circ}{J} \geq \overset{\circ}{J}$  case can be proven using the same technique.

**Proof of Lemma 3:**

We consider the mapping  $\mathbf{w}(\boldsymbol{\theta}) : \mathbb{R}^{\mathring{J}} \mapsto \mathbb{R}^J$ , where  $\mathring{J} < J$ , as an example.

Assume:

$\mathbf{w}(\boldsymbol{\theta})$  is continuously differentiable in  $\mathcal{N}_r(\boldsymbol{\theta}_0)$ .

If:

$$a_n(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0) \xrightarrow{d} V$$

where  $\boldsymbol{\theta}$  and  $V$  are some random vectors in  $\mathbb{R}^{\mathring{J}}$ ,  $\boldsymbol{\theta}_0$  is fixed in  $\mathbb{R}^{\mathring{J}}$ , and  $a_n \rightarrow \infty$ .

then:

$$a_n(\mathbf{w}(\hat{\boldsymbol{\theta}}) - \mathbf{w}(\boldsymbol{\theta}_0)) \xrightarrow{d} \nabla \mathbf{w}(\boldsymbol{\theta}) V$$

**Proof of Lemma 3:**

$\forall k \in \{1, \dots, J\}$ , extend  $\hat{\boldsymbol{\theta}}$  around  $\boldsymbol{\theta}_0$ , we have:

$$\left\{ \mathbf{w}_k(\hat{\boldsymbol{\theta}}) \right\}_{1 \times 1} = \left\{ \mathbf{w}_k(\boldsymbol{\theta}_0) \right\}_{1 \times 1} + \left\{ \nabla^{\boldsymbol{\theta}} \mathbf{w}_k(\boldsymbol{\theta}_k^*) \right\}_{1 \times \mathring{J}} \left\{ (\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0) \right\}_{\mathring{J} \times 1}$$

where  $\boldsymbol{\theta}_k^* \in \overline{(\boldsymbol{\theta}_0, \hat{\boldsymbol{\theta}})}$ , and  $\overline{(\boldsymbol{\theta}_0, \hat{\boldsymbol{\theta}})}$  denotes the open ball bounded by  $\boldsymbol{\theta}_0$  and  $\hat{\boldsymbol{\theta}}$ .

Define  $\left\{ \nabla^{\boldsymbol{\theta}} \mathbf{w} \right\}_{J \times \mathring{J}}$  as the partial derivative of  $\mathbf{w}_i$  with respect to  $\boldsymbol{\theta}$ .

Define  $P_n^* = \left\{ \nabla^{\boldsymbol{\theta}} \mathbf{w}_1(\boldsymbol{\theta}_1^*), \dots, \nabla^{\boldsymbol{\theta}} \mathbf{w}_J(\boldsymbol{\theta}_J^*) \right\}$ .

We have:

$$\mathbf{w}(\boldsymbol{\theta}) = \mathbf{w}(\boldsymbol{\theta}_0) + P_n^*(\boldsymbol{\theta} - \boldsymbol{\theta}_0).$$

Hence,  $\forall c \in \mathbb{R}^{J_1}$ :

$$c^\top a_n(\mathbf{w}(\hat{\boldsymbol{\theta}}) - \mathbf{w}(\boldsymbol{\theta}_0)) = c^\top (P_n^* - \nabla^{\boldsymbol{\theta}} \mathbf{w}(\boldsymbol{\theta}_0)) (a_n(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0)) + c^\top \nabla^{\boldsymbol{\theta}} \mathbf{w}(\boldsymbol{\theta}_0) a_n(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0).$$



While the first term converges to 0 in probability, the second term converges to  $c^\top \nabla^\theta \mathbf{w}(\boldsymbol{\theta}_0) V$  in distribution, by Cramér-Wold Theorem.

Therefore,  $\forall c \in \mathbb{R}^{J_1}$ , where  $J_1 > \overset{\circ}{J}$ :

$$c^\top a_n(\mathbf{w}(\boldsymbol{\theta}) - \mathbf{w}(\boldsymbol{\theta}_0)) \xrightarrow{d} c^\top \nabla^\theta \mathbf{w}(\boldsymbol{\theta}_0) V.$$

The result follows by applying the Cramér-Wold Theorem again. ■

## Proof of Theorem 2:

### (i. Consistency)

Under Regularity Condition (N-0), it can be shown:

$$\hat{\mathbf{w}}|\boldsymbol{\theta}, \lambda = (\lambda \mathbf{I} + \boldsymbol{\psi}(\boldsymbol{\theta}))^{-1} \boldsymbol{\phi}(\boldsymbol{\theta}) \quad (15)$$

$$\hat{\lambda}|\boldsymbol{\theta} = \arg_{\lambda} \left\{ [(\lambda \mathbf{I} + \boldsymbol{\psi}(\boldsymbol{\theta}))^{-1} \boldsymbol{\phi}(\boldsymbol{\theta})]^\top [(\lambda \mathbf{I} + \boldsymbol{\psi}(\boldsymbol{\theta}))^{-1} \boldsymbol{\phi}(\boldsymbol{\theta})] = 1 \right\} \quad (16)$$

$$\hat{\boldsymbol{\theta}}|\hat{\mathbf{w}}, \hat{\lambda} = \arg\max_{\boldsymbol{\theta}} L(\mathbf{D}; \hat{\mathbf{w}}, \boldsymbol{\theta}, \hat{\lambda}) \quad (17)$$

By (15) and (16), we have:

$$\hat{\mathbf{w}}|\boldsymbol{\theta} = \hat{\mathbf{w}}(\boldsymbol{\theta}). \quad (18)$$

Hence, for every fixed  $\hat{\boldsymbol{\theta}}$  inputting to (18), there is an unique output of  $\mathbf{w} : \mathbb{R}^{\overset{\circ}{J}} \mapsto \mathbb{R}^J$ , associated with that particular  $\mathbf{w}$ .

Specifically, for every fixed  $\mathbf{w}$  in (6),  $\exists \mathbf{w}$ , such that

$$\mathbf{w} = \mathbf{w}(\hat{\boldsymbol{\theta}}). \quad (19)$$

By the consistency proof of Theorem 1, under Regularity Condition (N-12), that  $\mathbf{w}(\boldsymbol{\theta})$  is continuously differentiable in  $\mathcal{N}_r(\boldsymbol{\theta}_0)$ , we have:

$$\mathbf{w}(\hat{\boldsymbol{\theta}}) := \hat{\mathbf{w}} \xrightarrow{p} \mathbf{w}(\boldsymbol{\theta}_0) := \mathbf{w}_0 \quad (20)$$

$\forall i \in \{1, \dots, q_1\}$ , by the continuous mapping theorem. ■

### (ii. Asymptotic Normality)

Under Regularity Condition (N-12) and (N-13), and by Lemma 3,

$$\sqrt{n}(\mathbf{w}(\hat{\boldsymbol{\theta}}) - \mathbf{w}(\boldsymbol{\theta}_0)) \rightarrow N\left(\nabla \mathbf{w}(\boldsymbol{\theta})\mu(\boldsymbol{\theta}_0), [\nabla \mathbf{w}(\boldsymbol{\theta})]^\top \Sigma(\boldsymbol{\theta}_0) \nabla \mathbf{w}(\boldsymbol{\theta})\right). \quad (21)$$

■

## 1.1 Obtain Variance Estimates of the first PDM

Recall for the first Principal Direction of Mediation, the lagrangian can be written as:

$$G(\mathbf{D}; \mathbf{w}, \boldsymbol{\theta}, \lambda) = g_1(\mathbf{D}; \mathbf{w}, \boldsymbol{\theta}) - \lambda(\mathbf{w}^\top \mathbf{w} - 1), \quad (22)$$

where  $g_1$  is the joint log-likelihood function.

Conditioning on  $\hat{\lambda}$  and  $\hat{\mathbf{w}}$ , the estimates of the multiplier and the first PDM, we have:

$$G(\mathbf{D}, \hat{\mathbf{w}}, \hat{\lambda}; \boldsymbol{\theta}) = g_1(\mathbf{D}, \hat{\mathbf{w}}; \boldsymbol{\theta}). \quad (23)$$

Define  $\dot{G}(\mathbf{D}, \hat{\mathbf{w}}, \hat{\lambda}; \boldsymbol{\theta}) = \frac{\partial g_1(\mathbf{D}, \hat{\mathbf{w}}; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}} = \sum_{i=1}^n \frac{\partial f_i(\mathbf{D}_i, \hat{\mathbf{w}}; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}}$ , where  $f_i(\mathbf{D}_i, \hat{\mathbf{w}}; \boldsymbol{\theta})$  is the log-likelihood function for  $\mathbf{D}_i = (X_i, Y_i, \mathbf{M}_i)$  conditioning on  $\mathbf{w} = \hat{\mathbf{w}}$ , i.e.,

$$f_i(\mathbf{D}_i, \hat{\mathbf{w}}; \boldsymbol{\theta}) = \frac{1}{2} \left\{ \frac{(Y_i - \beta_0 - X_i\gamma_1 - \mathbf{M}_i\hat{\mathbf{w}}\beta_1)^2}{\sigma_\epsilon} + \frac{(\mathbf{M}_i\hat{\mathbf{w}} - \alpha_0 - X_i\alpha_1)^2}{\sigma_\eta} \right\}.$$

The normalizing term  $\log \frac{1}{2\pi\sigma_\epsilon\sigma_\eta}$  and the minus sign in the log-likelihood function are dropped because they do not affect the final result.

$$\text{Define } \ell(D, \hat{\mathbf{w}}; \boldsymbol{\theta}) = \frac{\partial f(D, \hat{\mathbf{w}}; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}}.$$

### 1.1.1 Estimation of $D_0(\theta_0)$

Recall  $D_0(\theta_0) = \mathbb{E}_{\theta_0} \left( \frac{\partial \ell(D; \mathbf{w}, \theta)}{\partial \theta} \right)$ .

It is easy to see:

$$\ell(\mathbf{D}_i, \hat{\mathbf{w}}; \theta) = \frac{\partial f_i(\mathbf{D}_i, \hat{\mathbf{w}}; \theta)}{\partial \theta} = \begin{pmatrix} \frac{\partial f_i(\mathbf{D}_i, \hat{\mathbf{w}}; \theta)}{\partial \alpha_0} \\ \frac{\partial f_i(\mathbf{D}_i, \hat{\mathbf{w}}; \theta)}{\partial \alpha_1} \\ \frac{\partial f_i(\mathbf{D}_i, \hat{\mathbf{w}}; \theta)}{\partial \beta_0} \\ \frac{\partial f_i(\mathbf{D}_i, \hat{\mathbf{w}}; \theta)}{\partial \beta_1} \\ \frac{\partial f_i(\mathbf{D}_i, \hat{\mathbf{w}}; \theta)}{\partial \gamma_1} \end{pmatrix} = \begin{pmatrix} -\frac{(\mathbf{M}_i \mathbf{w} - \alpha_0 - X_i \alpha_1)}{\sigma_\eta^2} \\ -\frac{(\mathbf{M}_i \mathbf{w} - \alpha_0 - X_i \alpha_1)}{\sigma_\eta^2} X_i \\ -\frac{(Y_i - \beta_0 - X_i \gamma_1 - \mathbf{M}_i \mathbf{w} \beta_1)}{\sigma_\epsilon^2} \\ -\frac{(Y_i - \beta_0 - X_i \gamma_1 - \mathbf{M}_i \mathbf{w} \beta_1)}{\sigma_\epsilon^2} \mathbf{M}_i \mathbf{w} \\ -\frac{(Y_i - \beta_0 - X_i \gamma_1 - \mathbf{M}_i \mathbf{w} \beta_1)}{\sigma_\epsilon^2} X_i \end{pmatrix}. \quad (24)$$

Moreover,

$$\frac{\partial^2 \ell(\mathbf{D}_i, \hat{\mathbf{w}}; \theta)}{\partial \theta \partial \theta} = \begin{pmatrix} \frac{\partial^2 f_i}{\partial \alpha_0 \partial \alpha_0} & \frac{\partial^2 f_i}{\partial \alpha_0 \partial \alpha_1} & \frac{\partial^2 f_i}{\partial \alpha_0 \partial \beta_0} & \frac{\partial^2 f_i}{\partial \alpha_0 \partial \beta_1} & \frac{\partial^2 f_i}{\partial \alpha_0 \partial \gamma_1} \\ \frac{\partial^2 f_i}{\partial \alpha_1 \partial \alpha_0} & \frac{\partial^2 f_i}{\partial \alpha_1 \partial \alpha_1} & \frac{\partial^2 f_i}{\partial \alpha_1 \partial \beta_0} & \frac{\partial^2 f_i}{\partial \alpha_1 \partial \beta_1} & \frac{\partial^2 f_i}{\partial \alpha_1 \partial \gamma_1} \\ \frac{\partial^2 f_i}{\partial \beta_0 \partial \alpha_0} & \frac{\partial^2 f_i}{\partial \beta_0 \partial \alpha_1} & \frac{\partial^2 f_i}{\partial \beta_0 \partial \beta_0} & \frac{\partial^2 f_i}{\partial \beta_0 \partial \beta_1} & \frac{\partial^2 f_i}{\partial \beta_0 \partial \gamma_1} \\ \frac{\partial^2 f_i}{\partial \beta_1 \partial \alpha_0} & \frac{\partial^2 f_i}{\partial \beta_1 \partial \alpha_1} & \frac{\partial^2 f_i}{\partial \beta_1 \partial \beta_0} & \frac{\partial^2 f_i}{\partial \beta_1 \partial \beta_1} & \frac{\partial^2 f_i}{\partial \beta_1 \partial \gamma_1} \\ \frac{\partial^2 f_i}{\partial \gamma_1 \partial \alpha_0} & \frac{\partial^2 f_i}{\partial \gamma_1 \partial \alpha_1} & \frac{\partial^2 f_i}{\partial \gamma_1 \partial \beta_0} & \frac{\partial^2 f_i}{\partial \gamma_1 \partial \beta_1} & \frac{\partial^2 f_i}{\partial \gamma_1 \partial \gamma_1} \end{pmatrix} = \begin{pmatrix} \frac{1}{\sigma_\eta^2} & \frac{X_i}{\sigma_\eta^2} & 0 & 0 & 0 \\ \frac{X_i}{\sigma_\eta^2} & \frac{X_i^2}{\sigma_\eta^2} & 0 & 0 & 0 \\ 0 & 0 & \frac{1}{\sigma_\epsilon^2} & \frac{\mathbf{M}_i \mathbf{w}}{\sigma_\epsilon^2} & \frac{X_i}{\sigma_\epsilon^2} \\ 0 & 0 & \frac{\mathbf{M}_i \mathbf{w}}{\sigma_\epsilon^2} & \frac{(\mathbf{M}_i \mathbf{w})^2}{\sigma_\epsilon^2} & \frac{X_i \mathbf{M}_i \mathbf{w}}{\sigma_\epsilon^2} \\ 0 & 0 & \frac{X_i}{\sigma_\epsilon^2} & \frac{X_i \mathbf{M}_i \mathbf{w}}{\sigma_\epsilon^2} & \frac{X_i^2}{\sigma_\epsilon^2} \end{pmatrix}. \quad (25)$$

Therefore,  $D_0(\boldsymbol{\theta}_0) = \mathbb{E}_{\boldsymbol{\theta}_0} \left( \frac{\partial \ell(D; \mathbf{w}, \boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \right)$  can be estimated by

$$\hat{D}_0 = \begin{pmatrix} \frac{1}{\hat{\sigma}_\eta^2} & \frac{\bar{X}}{\hat{\sigma}_\eta^2} & 0 & 0 & 0 \\ \frac{\bar{X}}{\hat{\sigma}_\eta^2} & \frac{\sum_{i=1}^n X_i^2}{n} / \hat{\sigma}_\eta^2 & 0 & 0 & 0 \\ 0 & 0 & \frac{1}{\hat{\sigma}_\epsilon^2} & \frac{\sum_{i=1}^n \mathbf{M}_i \hat{\mathbf{w}}}{n} / \hat{\sigma}_\epsilon^2 & \frac{\bar{X}}{\hat{\sigma}_\epsilon^2} \\ 0 & 0 & \frac{\sum_{i=1}^n \mathbf{M}_i \hat{\mathbf{w}}}{n} / \hat{\sigma}_\epsilon^2 & \frac{\sum_{i=1}^n (\mathbf{M}_i \hat{\mathbf{w}})^2}{n} / \hat{\sigma}_\epsilon^2 & \frac{\sum_{i=1}^n X_i \mathbf{M}_i \hat{\mathbf{w}}}{n} / \hat{\sigma}_\epsilon^2 \\ 0 & 0 & \frac{\bar{X}}{\hat{\sigma}_\epsilon^2} & \frac{\sum_{i=1}^n X_i \mathbf{M}_i \hat{\mathbf{w}}}{n} / \hat{\sigma}_\epsilon^2 & \frac{\sum_{i=1}^n X_i^2}{n} / \hat{\sigma}_\epsilon^2 \end{pmatrix}. \quad (26)$$

### 1.1.2 Estimation of $V(\boldsymbol{\theta}_0)$

Recall  $V(\boldsymbol{\theta}_0) = \mathbb{E}_{\boldsymbol{\theta}_0} \left( \ell(D; \mathbf{w}, \boldsymbol{\theta}) \ell^\top(D; \mathbf{w}, \boldsymbol{\theta}) \right) - m(\boldsymbol{\theta}_0) m^\top(\boldsymbol{\theta}_0)$ .

For simplicity, denote  $f_i = f_i(\mathbf{D}_i, \hat{\mathbf{w}}; \boldsymbol{\theta})$ . It is easy to see:

$$\ell(\mathbf{D}_i, \hat{\mathbf{w}}; \boldsymbol{\theta}) \ell^\top(\mathbf{D}_i, \hat{\mathbf{w}}; \boldsymbol{\theta}) = \begin{pmatrix} \left[ \frac{\partial f_i}{\partial \alpha_0} \right]^2 & \frac{\partial f_i}{\partial \alpha_0} \frac{\partial f_i}{\partial \alpha_1} & \frac{\partial f_i}{\partial \alpha_0} \frac{\partial f_i}{\partial \beta_0} & \frac{\partial f_i}{\partial \alpha_0} \frac{\partial f_i}{\partial \beta_1} & \frac{\partial f_i}{\partial \alpha_0} \frac{\partial f_i}{\partial \gamma_1} \\ \frac{\partial f_i}{\partial \alpha_1} \frac{\partial f_i}{\partial \alpha_0} & \left[ \frac{\partial f_i}{\partial \alpha_1} \right]^2 & \frac{\partial f_i}{\partial \alpha_1} \frac{\partial f_i}{\partial \beta_0} & \frac{\partial f_i}{\partial \alpha_1} \frac{\partial f_i}{\partial \beta_1} & \frac{\partial f_i}{\partial \alpha_1} \frac{\partial f_i}{\partial \gamma_1} \\ \frac{\partial f_i}{\partial \beta_0} \frac{\partial f_i}{\partial \alpha_0} & \frac{\partial f_i}{\partial \beta_0} \frac{\partial f_i}{\partial \alpha_1} & \left[ \frac{\partial f_i}{\partial \beta_0} \right]^2 & \frac{\partial f_i}{\partial \beta_0} \frac{\partial f_i}{\partial \beta_1} & \frac{\partial f_i}{\partial \beta_0} \frac{\partial f_i}{\partial \gamma_1} \\ \frac{\partial f_i}{\partial \beta_1} \frac{\partial f_i}{\partial \alpha_0} & \frac{\partial f_i}{\partial \beta_1} \frac{\partial f_i}{\partial \alpha_1} & \frac{\partial f_i}{\partial \beta_1} \frac{\partial f_i}{\partial \beta_0} & \left[ \frac{\partial f_i}{\partial \beta_1} \right]^2 & \frac{\partial f_i}{\partial \beta_1} \frac{\partial f_i}{\partial \gamma_1} \\ \frac{\partial f_i}{\partial \gamma_1} \frac{\partial f_i}{\partial \alpha_0} & \frac{\partial f_i}{\partial \gamma_1} \frac{\partial f_i}{\partial \alpha_1} & \frac{\partial f_i}{\partial \gamma_1} \frac{\partial f_i}{\partial \beta_0} & \frac{\partial f_i}{\partial \gamma_1} \frac{\partial f_i}{\partial \beta_1} & \left[ \frac{\partial f_i}{\partial \gamma_1} \right]^2 \end{pmatrix} \quad (27)$$

$$= \begin{pmatrix} R_i^2 & R_i^2 X_i & R_i S_i & R_i S_i \mathbf{M}_i \mathbf{w} & R_i S_i X_i \\ R_i^2 X_i & R_i^2 X_i^2 & R_i S_i X_i & R_i S_i \mathbf{M}_i \mathbf{w} X_i & R_i S_i X_i^2 \\ R_i S_i & R_i S_i X_i & S_i^2 & S_i^2 \mathbf{M}_i \mathbf{w} & S_i^2 X_i \\ R_i S_i \mathbf{M}_i \mathbf{w} & R_i S_i \mathbf{M}_i \mathbf{w} X_i & S_i \mathbf{M}_i \mathbf{w} & S_i^2 (\mathbf{M}_i \mathbf{w})^2 & S_i^2 \mathbf{M}_i \mathbf{w} X_i \\ R_i S_i X_i & R_i S_i X_i^2 & S_i^2 X_i & S_i^2 \mathbf{M}_i \mathbf{w} X_i & S_i^2 X_i^2 \end{pmatrix}$$

where  $R_i = -\frac{(\mathbf{M}_i \mathbf{w} - \alpha_0 - X_i \alpha_1)}{\sigma_\eta^2}$  and  $S_i = -\frac{(Y_i - \beta_0 - X_i \gamma_1 - \mathbf{M}_i \mathbf{w} \beta_1)}{\sigma_\epsilon^2}$ .

Since  $m(\boldsymbol{\theta}_0) = O_p(1)$ ,  $V(\boldsymbol{\theta}_0)$  can be estimated by

$$\hat{V}(\hat{\boldsymbol{\theta}}) = \begin{pmatrix} \frac{\sum_{i=1}^n \hat{R}_i^2}{n} & \frac{\sum_{i=1}^n \hat{R}_i^2 X_i}{n} & \frac{\sum_{i=1}^n \hat{R}_i \hat{S}_i}{n} & \frac{\sum_{i=1}^n \hat{R}_i \hat{S}_i \mathbf{M}_i \mathbf{w}}{n} & \frac{\sum_{i=1}^n \hat{R}_i \hat{S}_i X_i}{n} \\ \frac{\sum_{i=1}^n \hat{R}_i^2 X_i}{n} & \frac{\sum_{i=1}^n \hat{R}_i^2 X_i^2}{n} & \frac{\sum_{i=1}^n \hat{R}_i \hat{S}_i X_i}{n} & \frac{\sum_{i=1}^n \hat{R}_i \hat{S}_i \mathbf{M}_i \hat{\mathbf{w}} X_i}{n} & \frac{\sum_{i=1}^n \hat{R}_i \hat{S}_i X_i^2}{n} \\ \frac{\sum_{i=1}^n \hat{R}_i \hat{S}_i}{n} & \frac{\sum_{i=1}^n \hat{R}_i \hat{S}_i X_i}{n} & \frac{\sum_{i=1}^n \hat{S}_i^2}{n} & \frac{\sum_{i=1}^n \hat{S}_i^2 \mathbf{M}_i \hat{\mathbf{w}}}{n} & \frac{\sum_{i=1}^n \hat{S}_i^2 X_i}{n} \\ \frac{\sum_{i=1}^n \hat{R}_i \hat{S}_i \mathbf{M}_i \hat{\mathbf{w}}}{n} & \frac{\sum_{i=1}^n \hat{R}_i \hat{S}_i \mathbf{M}_i \hat{\mathbf{w}} X_i}{n} & \frac{\sum_{i=1}^n \hat{S}_i \mathbf{M}_i \hat{\mathbf{w}}}{n} & \frac{\sum_{i=1}^n \hat{S}_i^2 (\mathbf{M}_i \hat{\mathbf{w}})^2}{n} & \frac{\sum_{i=1}^n \hat{S}_i^2 \mathbf{M}_i \hat{\mathbf{w}} X_i}{n} \\ \frac{\sum_{i=1}^n \hat{R}_i \hat{S}_i X_i}{n} & \frac{\sum_{i=1}^n \hat{R}_i \hat{S}_i X_i^2}{n} & \frac{\sum_{i=1}^n \hat{S}_i^2 X_i}{n} & \frac{\sum_{i=1}^n \hat{S}_i^2 \mathbf{M}_i \hat{\mathbf{w}} X_i}{n} & \frac{\sum_{i=1}^n \hat{S}_i^2 X_i^2}{n} \end{pmatrix}. \quad (28)$$

where  $\hat{\mathbf{w}}$  is the estimates of the first PDM, and  $\hat{R}_i = -\frac{(\mathbf{M}_i \hat{\mathbf{w}} - \hat{\alpha}_0 - X_i \hat{\alpha}_1)}{\hat{\sigma}_\eta^2}$  and  $\hat{S}_i =$

$$-\frac{(Y_i - \hat{\beta}_0 - X_i \hat{\gamma}_1 - \mathbf{M}_i \hat{\mathbf{w}} \hat{\beta}_1)}{\hat{\sigma}_\epsilon^2}.$$

### 1.1.3 Estimation of Asymptotic Variance for Pathway Coefficients

Hence, we can estimate the asymptotic variance for pathway coefficients,  $\Sigma(\boldsymbol{\theta}_0)$ , by

$$\hat{\Sigma}(\hat{\boldsymbol{\theta}}) = \hat{D}_0^{-1} \hat{V}(\hat{\boldsymbol{\theta}}) [\hat{D}_0^{-1}]^\top \quad (29)$$

where  $\hat{D}_0$  and  $\hat{V}(\hat{\boldsymbol{\theta}})$  are based upon (26) and (28).

### 1.1.4 Estimation of Asymptotic Variance for the first PDM

First, we need to find the first partial derivative of  $\mathbf{w}$  with respect to  $\boldsymbol{\theta}$ .

Due to (11) in the paper, we have:

$$\mathbf{w}|\boldsymbol{\theta}, \hat{\lambda} = (\hat{\lambda} \mathbf{I} + \psi(\boldsymbol{\theta}))^{-1} \phi(\boldsymbol{\theta}) \quad (30)$$

where

$$\psi(\boldsymbol{\theta}) = \frac{\mathbf{M}^\top \mathbf{M} \beta_1^2}{\sigma_{\epsilon_1}^2} + \frac{\mathbf{M}^\top \mathbf{M}}{\sigma_{\eta_1}^2} \quad (31)$$

and

$$\phi(\boldsymbol{\theta}) = \frac{\mathbf{M}^\top (\alpha_{0,1} + \alpha_1 \mathbf{X})}{\sigma_{\eta_1}^2} + \frac{\mathbf{M}^\top (\mathbf{Y} - \beta_{0,1} - \mathbf{X} \gamma_1) \beta_1}{\sigma_{\epsilon_1}^2}. \quad (32)$$

.

**Lemma 4** For a non-singular square matrix  $\mathbf{A}$ , where elements are a function of  $\boldsymbol{\theta}$ , then:

$$\frac{\partial \mathbf{A}^{-1}}{\partial \boldsymbol{\theta}} = -\mathbf{A}^{-1} \frac{\partial \mathbf{A}}{\partial \boldsymbol{\theta}} \mathbf{A}^{-1}.$$

Due to Lemma 4, we have  $\nabla^{\mathbf{w}}(\boldsymbol{\theta}) = \frac{\partial \mathbf{w}}{\partial \boldsymbol{\theta}} = -\mathbf{A}^{-1} \frac{\partial \psi(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \mathbf{A}^{-1} \phi(\boldsymbol{\theta}) + \mathbf{A}^{-1} \frac{\partial \phi(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}}$ , where  $\mathbf{A} = \hat{\lambda} \mathbf{I} + \psi(\boldsymbol{\theta})$ .

With some algebra we can show that:

$$\nabla^{\mathbf{w}}(\boldsymbol{\theta}) = \begin{pmatrix} \frac{\partial \mathbf{w}}{\partial \alpha_0} \\ \frac{\partial \mathbf{w}}{\partial \alpha_1} \\ \frac{\partial \mathbf{w}}{\partial \beta_0} \\ \frac{\partial \mathbf{w}}{\partial \beta_1} \\ \frac{\partial \mathbf{w}}{\partial \gamma_1} \end{pmatrix} = \begin{pmatrix} \mathbf{A}^{-1} \frac{\mathbf{M}^\top \mathbf{1}}{\sigma_\eta^2} \\ \mathbf{A}^{-1} \frac{\mathbf{M}^\top \mathbf{X}}{\sigma_\eta^2} \\ \mathbf{A}^{-1} \left( -\frac{\mathbf{M}^\top \mathbf{1} \beta_1}{\sigma_\epsilon^2} \right) \\ -\mathbf{A}^{-1} \left[ \frac{2\beta_1 \mathbf{M}^\top \mathbf{M}}{\sigma_\epsilon^2} \right] \mathbf{A}^{-1} \phi(\boldsymbol{\theta}) + \mathbf{A}^{-1} \frac{\mathbf{M}^\top (\mathbf{Y} - \beta_0 - \mathbf{X} \gamma_1)}{\sigma_\epsilon^2} \\ \mathbf{A}^{-1} \left( -\frac{\mathbf{M}^\top \mathbf{X} \beta_1}{\sigma_\epsilon^2} \right) \end{pmatrix} \quad (33)$$

where  $\mathbf{1} = \text{vec}_{n \times 1}(1)$ .

Hence,  $\nabla^{\mathbf{w}}(\boldsymbol{\theta}_0)$  can be estimated by

$$\hat{\nabla}^{\mathbf{w}}(\hat{\boldsymbol{\theta}}) = \begin{pmatrix} \hat{\mathbf{A}}^{-1} \frac{\mathbf{M}^\top \mathbf{1}}{\hat{\sigma}_\eta^2} \\ \hat{\mathbf{A}}^{-1} \frac{\mathbf{M}^\top \mathbf{X}}{\hat{\sigma}_\eta^2} \\ \hat{\mathbf{A}}^{-1} \left( -\frac{\mathbf{M}^\top \mathbf{1} \hat{\beta}_1}{\hat{\sigma}_\epsilon^2} \right) \\ -\hat{\mathbf{A}}^{-1} \left[ \frac{2\hat{\beta}_1 \mathbf{M}^\top \mathbf{M}}{\hat{\sigma}_\epsilon^2} \right] \hat{\mathbf{A}}^{-1} \hat{\phi}(\hat{\boldsymbol{\theta}}) + \hat{\mathbf{A}}^{-1} \frac{\mathbf{M}^\top (\mathbf{Y} - \hat{\beta}_0 - \mathbf{X} \hat{\gamma}_1)}{\hat{\sigma}_\epsilon^2} \\ \hat{\mathbf{A}}^{-1} \left( -\frac{\mathbf{M}^\top \mathbf{X} \hat{\beta}_1}{\hat{\sigma}_\epsilon^2} \right) \end{pmatrix} \quad (34)$$

where  $\hat{\mathbf{A}} = \hat{\lambda} \mathbf{I} + \hat{\psi}(\hat{\boldsymbol{\theta}})$ , and  $\hat{\boldsymbol{\theta}} = \{\hat{\alpha}_0, \hat{\alpha}_1, \hat{\beta}_0, \hat{\beta}_1, \hat{\gamma}_1\}$ .

Therefore, the asymptotic variance for the first PDM,  $\Sigma^{\mathbf{w}}(\boldsymbol{\theta}_0)$  can be estimated by:

$$\hat{\Sigma}^{\mathbf{w}}(\hat{\boldsymbol{\theta}}) = [\hat{\nabla}^{\mathbf{w}}(\hat{\boldsymbol{\theta}})]^\top \hat{D}_0^{-1} \hat{V}(\hat{\boldsymbol{\theta}}) [\hat{D}_0^{-1}]^\top \hat{\nabla}^{\mathbf{w}}(\hat{\boldsymbol{\theta}}) \quad (35)$$

with each term on the right side based upon (26), (28), and (34).

### 1.1.5 Estimation of the Conditional Asymptotic Variance for the first PDM in a Ultra-high Dimensional Setting

When the mediation matrix  $\mathbf{M}$  is ultrahigh dimensional, we first conduct Population Value Decomposition (PVD), as explained in the paper. We obtain  $\tilde{\mathbf{M}}$  and  $\tilde{\mathbf{w}}$ , where  $\mathbf{M}\mathbf{w} = \tilde{\mathbf{M}}\tilde{\mathbf{w}}$  and  $\tilde{\mathbf{w}} = \mathbf{D}\mathbf{w}$ , where  $\mathbf{D}$  is the matrix carrying population level information, obtained via the PVD in the paper. We then obtain the asymptotic variance for  $\tilde{\mathbf{w}}$  using the aforementioned method, say,  $\hat{\Sigma}^{\tilde{\mathbf{w}}}(\hat{\boldsymbol{\theta}})$ , and it follows that estimate of the conditional asymptotic variance for  $\mathbf{w}$  is

$$\hat{\Sigma}^{\mathbf{w}}(\hat{\boldsymbol{\theta}}|\mathbf{D}) = \mathbf{D}^{-1} \frac{\hat{\Sigma}^{\tilde{\mathbf{w}}}(\hat{\boldsymbol{\theta}})}{n} [\mathbf{D}^{\top}]^{-1}. \quad (36)$$

It is necessary to mention that,  $\mathbf{D}$  obtained from PVD is random. If it were not conditioned, it is difficulty to obtain  $\hat{\Sigma}^{\mathbf{w}}(\hat{\boldsymbol{\theta}})$  in that when  $\mathbf{D}$  is of ultrahigh dimensional, it is computationally costly to estimate the variance-covariance matrix of  $\mathbf{D}$ .

## References and Notes

Newey, W. K. and D. McFadden (1994). Large sample estimation and hypothesis testing. *Handbook of econometrics* 4, 2111–2245.