

Time-varying mediation analysis for incomplete data with application to DNA methylation study for PTSD

Kecheng Wei^a, Fei Xue^b, Qi Xu^c, Yubai Yuan^d, Yuexia Zhang^e,
Guoyou Qin^a, Agaz H. Wani^f, Allison E. Aiello^g,
Derek E. Wildman^f, Monica Uddin^f and Annie Qu^{h*}

^aDepartment of Biostatistics, Fudan University;

^bDepartment of Statistics, Purdue University;

^cDepartment of Statistics and Data Science, Carnegie Mellon University;

^dDepartment of Statistics, Penn State University;

^eDepartment of Management Science and Statistics, The University of Texas at San Antonio;

^fCollege of Public Health, University of South Florida;

^gDepartment of Epidemiology and Robert N. Butler Columbia Aging Center, Columbia University

^hDepartment of Statistics, University of California Irvine;

Abstract

DNA methylation (DNAm) has been shown to mediate causal effects from traumatic experiences to post-traumatic stress disorder (PTSD). However, the scientific question about whether the mediation effect changes over time remains unclear. In this paper, we develop time-varying structural equation models to identify cytosine-phosphate-guanine (CpG) sites where DNAm mediates the effect of trauma exposure on PTSD, and to capture dynamic changes in mediation effects. The proposed methodology is motivated by the Detroit Neighborhood Health Study (DNHS) with high-dimensional and longitudinal DNAm measurements. To handle the non-monotone missing DNAm in the dataset, we propose a novel Longitudinal Multiple Imputation (LMI) method utilizing dependency among repeated measurements, and employ the generalized method of moments to integrate the multiple imputations. Simulations confirm that the proposed method outperforms existing approaches in various longitudinal settings. In DNHS data analysis, our method identifies several CpG sites where DNAm exhibits dynamic mediation effects. Some of the corresponding genes have been shown to be associated with PTSD in the existing literature, and our findings

*CONTACT Annie Qu. Email: aqu2@uci.edu

on their time-varying effects could deepen the understanding of the mediation role of DNAm on the causal path from trauma exposure to PTSD risk.

Keywords: generalized method of moments, high dimensional mediators, longitudinal data, missing data, multiple imputation, structural equation modeling, variable selection.

1 Introduction

Post-traumatic stress disorder (PTSD) is a serious mental health disorder that arises after an individual witnesses or experiences traumatic events, such as severe accidents, natural disasters, or warfare. Although the causality of trauma on PTSD has been established from the biological perspective, the trauma-PTSD relation has not been elucidated on the population level of public health data. World Health Organization World Mental Health Surveys show that, although 69.7% of the surveyed population were exposed to at least one traumatic event, the cross-national lifetime prevalence of PTSD was only 5.6% among the trauma-exposed subpopulation (Koenen et al. 2017). This significant disparity of prevalence between trauma exposure and PTSD implies that there are potential intermediate factors intervening in the trauma-PTSD pathway to explain individual vulnerability to trauma and tendency to develop PTSD. On the other hand, extensive meta-analyses show that the variability in PTSD on the population level is not sufficiently explained by incorporating phenotype variables as predictors (Tortella-Feliu et al. 2019). These limitations in previous PTSD studies lead us to further explore genotype information, such as genetic and relevant molecular variation, as intermediate variables, and identify causal mediation pathways from trauma exposure to PTSD via the transmission of genetic variation.

DNA methylation (DNAm) is a molecular modification that alters gene expression without changing DNA sequence, and can be influenced by stress and socio-environmental factors (Turecki and Meaney 2016, Lussier et al. 2023). It has been found that aberrant DNAm is associated with a variety of stress-related neuropsychiatric disorders such as depression, schizophrenia, and PTSD (Wani et al. 2021, Wilker et al. 2023). Thus, DNAm may be a pathway by which traumatic events become biologically embedded and contribute to mental illness. Given the intermediate role of DNAm in mental disorder formulation, studies have suggested that DNAm may also mediate the trauma-PTSD mechanism: trauma increases or reduces DNAm levels, which in turn exacerbate or alleviate PTSD symptoms.

For example, [Rutten et al. \(2018\)](#) found that DNAm mediates the relationship between combat trauma and PTSD. [Occean et al. \(2022\)](#) showed that DNAm in the nuclear factor of activated T cells 1 mediates the exposure to different types of traumatic events and the risk for PTSD. [Xue et al. \(2022\)](#) identified the heterogeneous mediation effects of DNAm in different subpopulations.

From the perspective of mediation methodology, most of the existing mediation analyses on DNAm are cross-sectional based, i.e., only estimating mediation pathways at a single time point. However, DNAm levels and PTSD severity typically progress over time, therefore the relationships among trauma, DNAm, and PTSD may temporally change. For example, [Wu et al. \(2021\)](#) found that different pre-migration and post-migration stressors significantly affect refugees' mental health at different resettlement stages. [Lussier et al. \(2023\)](#) demonstrated the temporal effects of adversity exposure on genetic variation across childhood and adolescence. [Wilker et al. \(2023\)](#) showed that the association between glucocorticoid receptor gene methylation status and PTSD symptoms is only significant in the early stages of study. Nevertheless, dynamic mediation analysis tools are still not well developed. By performing time-varying mediation analysis on longitudinal data, we can achieve a more comprehensive understanding of how causal pathways between variables in the disease process evolve over time: when the effects occur, diminish, strengthen, weaken, or change direction. Such insights could enrich understanding of the dynamic mediation effects of DNAm on pathways from trauma exposure to PTSD risk, and help to develop more targeted and time-sensitive interventions.

There have been some pioneering works about time-varying mediation analysis where exposure, mediator, and outcome are observed over multiple time points. For example, [Rijnhart et al. \(2022\)](#) and [Cai et al. \(2022\)](#) included time-interaction terms in their model to account for temporal effects. [Zeng et al. \(2021\)](#) and [Zeng et al. \(2022\)](#) extended the classical causal mediation method from a functional data analysis perspective. [Ge et al. \(2023\)](#) and [Luo et al. \(2023\)](#) studied dynamic mediation effects under a reinforcement learning framework. Nevertheless, most focus on cases where there is only one or a small number of mediators, with completely observed data. Few studies consider scenarios of high-dimensional mediators with incomplete data.

Missing data is ubiquitous in longitudinal study. In the statistical literature, two types of missingness structures are generally considered. One is called “monotone missingness”

or “dropout”, where a subject may leave the study at a certain time point and never return. On the other hand, “non-monotone missingness” or “intermittent missingness” refers to when a subject may miss certain visits during follow-up and return at later visits. For our motivating Detroit Neighborhood Health Study (DNHS) which will be introduced later, subjects are supposed to have the methylation status of a large number of cytosine-phosphate-guanine (CpG) sites measured at each follow-up visit. Due to variability of compliance, the DNAm data exhibits the non-monotone missingness mechanism, and high-dimensional DNAm values are entirely missing at some visits by the individual. Classical approaches, such as weighting ([Chen et al. 2021](#)), imputation ([Jahangiri et al. 2023](#)), and likelihood-based methods ([Tseng et al. 2016](#)), have been developed to address non-monotone missingness in longitudinal data. However, these approaches typically consider non-monotone missingness of the outcome or with a few covariates, and may not be directly applicable to our setting where high-dimensional DNAm values are entirely unobserved at some visits.

In this paper, we develop a new time-varying mediation method to investigate the temporal mediation effects of DNAm transmission from trauma to PTSD. We propose time-varying structural equation models, and adopt a regularization approach to select relevant mediators and identify time-varying effects, which yields several advantages. First, we can identify the important mediators from high-dimensional DNAm values at different time points, in contrast to existing dynamic mediation approaches considering only one or a small number of mediators ([Rijnhart et al. 2022](#)). Second, we can capture the trajectory of mediation effects even with a small number of time points, while existing approaches require dense time grids ([Cai et al. 2022](#)).

To handle non-monotone missing DNAm values, we propose a new imputation method called Longitudinal Multiple Imputation (LMI), which has the following major advantages. First, it makes full use of the dependence among repeated measurements, which exploits the inherent correlation structure of longitudinal data, not just based on low-rank structures or trajectory means ([Mazumder, Hastie and Tibshirani 2010](#), [Jahangiri et al. 2023](#)). Second, it can handle high-dimensional mediators being entirely unobserved at some time points, in contrast to existing approaches which only consider missingness of the outcome or with a few covariates ([Chen et al. 2021](#)). Third, it does not require some individuals with complete repeated measurements in the dataset, whereas existing approaches may require that there

must be some subjects with complete data, or that the data at the first time visit is always observed for all subjects (Jahangiri et al. 2023).

In simulation studies, the proposed method outperforms existing approaches under different missingness structures and different dimensions of mediators, and its performance is robust to different missingness mechanisms. In DNHS data analysis, our method identifies some DNAm CpG sites which exhibit dynamic mediation effects. Specifically, certain CpG sites initially have nonzero mediation effects, but the effects disappear over time. Certain CpG sites do not show mediation effects in early stages, but the effects emerge over time. Some genes (e.g., *PTPRK*) where selected CpG sites (e.g., cg09247979) are located have been reported as differentially expressed genes in previous PTSD meta-analyses (Chitrala, Nagarkatti and Nagarkatti 2016), and our further findings on their time-varying effects elucidate trauma-PTSD dynamic causal mechanisms.

The rest of the paper is organized as follows. Section 2 describes the background of the DNHS dataset. Section 3 introduces the proposed time-varying mediation models and the non-monotone missingness structure of mediators. Section 4 presents the proposed approaches for multiple imputation, data integration, and regularization. Section 5 illustrates the detailed implementation and algorithm. Section 6 includes the results of simulation studies. Section 7 applies the proposed method to the DNHS dataset. Section 8 concludes the research with some discussion.

2 Motivating application: Longitudinal trauma, PTSD, and DNAm

2.1 DNHS data

The Detroit Neighborhood Health Study (DNHS) is a prospective and representative longitudinal cohort study of predominantly African American adults living in Detroit, Michigan (Uddin et al. 2010, Goldmann et al. 2011). Participants completed a 40-minute structured telephone interview annually between 2008-2012 (i.e., waves 1-5), to assess demographic characteristics, exposure to traumatic events, and mental health status. All participants were given the opportunity to provide a specimen (venipuncture, blood spot, or saliva) for immune and inflammatory marker testing as well as genetic testing of DNA. We analyze

the DNAm data with 526 subjects from waves 1, 2, 4, and 5 as there was no full blood draw in wave 3 as in the other waves. We refer to waves 1, 2, 4, and 5 as time points 1-4 throughout the paper.

The exposure variable of interest in our study is the number of traumatic events. In the initial interview (time point 1), participants were asked about their lifetime traumatic experiences using a predefined list of 19 traumatic events (Breslau et al. 1998). In subsequent interviews (time points 2-4), they were additionally queried about traumatic experiences since their last interview within the same list. At time points 2-4, we incorporate the number of traumatic events experienced since the previous interview to capture the cumulative severity of trauma exposure. The first plot in the left panel of Figure 1 shows the longitudinal trajectories of trauma from 15 subjects.

The outcome variable of interest is PTSD symptom severity, which is measured through the widely used self-report PTSD Checklist (PCL, Ruggiero et al. 2003). The PCL contains 17 items corresponding to key symptoms of PTSD. In the initial interview (time point 1), participants were asked how much they had been bothered by each symptom using a 5-point scale (1-5) in reference to their worst lifetime traumatic experience, and the summary score for the 17 items were used to characterize the PTSD symptom severity. In subsequent interviews (time points 2-4), they were additionally queried about the worst traumatic experience since their last interview, and the summary score corresponding to it was similarly calculated. At time points 2-4, we choose the largest score calculated in previous and current interviews to capture the cumulative severity of PTSD symptoms. The second plot in the left panel of Figure 1 shows the longitudinal trajectories of PTSD from 15 subjects. We additionally incorporate demographic information such as age and gender, and blood work such as CD4 T cells and B cells, measured at baseline.

Existing studies suggest that epigenomic variations can be induced by trauma exposure and accompany the development of stress-related disorders (Rutten et al. 2018, O'cean et al. 2022, Xue et al. 2022). However, few studies have reported brain-related epigenomic profiles that associate with PTSD risk. In DNHS, the methylation levels of 1879 CpG sites from paired blood and brain tissue (Braun et al. 2019) were assessed. We treat the DNAm levels of these 1879 CpG sites as potential mediators. The methylation level of a CpG site is quantified using the “Beta-value”, which ranges from 0 to 1 and is calculated as the proportion $M/(M+U+100)$, where M and U represent the signal intensities of methylated

and unmethylated probes (Mou et al. 2022). The third and fourth plots in the left panel of Figure 1 show the longitudinal trajectories of DNAm levels for two CpG sites from 15 subjects. We observe that the DNAm levels fluctuate within a specific range over the follow-up period, without exhibiting a clear trend or consistent pattern. The average rate of change in DNAm levels for all 1879 CpG sites across the 15 subjects between time points t and t' can be calculated as $CR_{tt'} = \frac{1}{15} \sum_{i=1}^{15} \sum_{j=1}^{1879} (M_{it}^j - M_{it'}^j) / M_{it'}^j$. The corresponding results for CR_{21} , CR_{31} , CR_{41} , CR_{32} , CR_{42} , and CR_{43} are 2.6%, 0%, 1.9%, -0.1%, 1.2%, and 3.2%, respectively.

The traumatic events and PTSD severity are easily assessed, so that we have complete repeated measurements for exposure and outcome. However, collecting and testing biospecimens is difficult and costly. Due to variability of compliance, loss to follow-up, and budget restriction, DNAm data could suffer from missingness. In DNHS, DNAm data of some subjects may be available only at certain visits, but are unobserved at other visits. The right panel of Figure 1 shows the non-monotone missingness structure of DNAm. Each blank column represents unobserved DNAm, while the colored ones represent observed DNAm. Participants can be divided into 15 groups, where group 1 has 4 repeated observations at all time points, while group 2 only has observations at time points 1-3, with high-dimensional DNAm data entirely missing at time point 4. The sample sizes of groups 1-15 are (15, 29, 7, 28, 21, 45, 27, 12, 23, 9, 21, 147, 24, 102, 16), respectively. There are 511 out of 526 subjects with missing values and the missingness rate of DNAm is 58.3%.

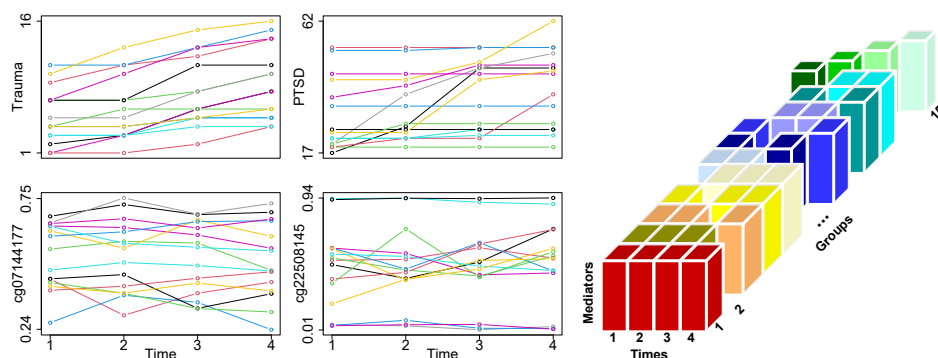


Figure 1: Left: Longitudinal trajectories of trauma, PTSD, and DNAm of two CpG sites from 15 subjects with complete repeated measurements of DNAm. Right: Non-monotone missingness structure of DNAm across all time points. Each blank column represents unobserved DNAm, while the colored ones represent observed DNAm.

2.2 Novelty of our analysis

Several studies have suggested that DNAm plays a critical role as a mediator in the causal relationship between trauma and PTSD. For instance, [Xue et al. \(2022\)](#) leveraged DNHS data to identify multiple DNAm CpG sites mediating the effect of trauma on PTSD, highlighting that these epigenetic factors can play different mediation roles across different sub-populations. However, their analysis relied on baseline wave to estimate static mediation effects, which is not applicable for discovering potential time-varying genetic contributions to the disease process.

Time-varying effects are a common phenomenon in the genetic regulation of disease progression. Dynamic factors such as environmental changes and the natural course of the disease interact intricately with gene regulation mechanisms, suggesting that genetic variants may adapt to changing environmental conditions and show different effects at different stages of the disease ([Lussier et al. 2023](#)). Investigating time-varying genetic effects provides a more nuanced understanding of how genetic factors influence disease processes, helps identify critical time windows for intervention, and facilitates the design of time-sensitive therapeutic strategies.

However, studying time-varying effects poses great challenges. As illustrated in Figure 1, while trauma exposure and PTSD severity tend to increase monotonically over time, the dynamics of DNAm appear irregular, making it difficult to capture their patterns of variation. Moreover, the issues of missingness and high dimensionality within the DNHS dataset further complicate the analysis. These challenges underscore the necessity of developing innovative statistical models and estimation methods to characterize time-varying effects and accommodate complex data structures.

3 Models

In this section, we propose time-varying structural equation models to study dynamic mediation effects, and introduce the non-monotone missingness structure of mediators.

3.1 Time-varying structural equation models

Let $\mathbf{M}_{it} = (M_{it}^1, \dots, M_{it}^p)^T$ be a $p \times 1$ vector of mediators, where M_{it}^j is the j -th ($j = 1, \dots, p$) mediator of the i -th ($i = 1, \dots, N$) subject, measured at the t -th ($t = 1, \dots, T$) time point.

Let $X_{it} \in \mathbb{R}$, $Y_{it} \in \mathbb{R}$, and $\mathbf{Z}_i \in \mathbb{R}^q$ be the longitudinal exposure, the longitudinal outcome, and the demographic variables measured at baseline, respectively. We propose the following structural equation models:

$$Y_{it} = \alpha_t X_{it} + \beta_t^T \mathbf{M}_{it} + \delta^T \mathbf{Z}_i + e_{it}, \quad (1)$$

$$\mathbf{M}_{it} = \gamma_t X_{it} + \boldsymbol{\eta}_t Y_{i,t-1} + \mathbf{\Gamma} \mathbf{Z}_i + \boldsymbol{\varepsilon}_{it}, \quad (2)$$

with $Y_{i0} = 0$, where the coefficients α_t , $\beta_t = (\beta_t^1, \dots, \beta_t^p)^T$, $\gamma_t = (\gamma_t^1, \dots, \gamma_t^p)^T$, and $\boldsymbol{\eta}_t = (\eta_t^1, \dots, \eta_t^p)^T$ measure the time-varying effects at the t -th time point, $\delta \in \mathbb{R}^q$ and $\mathbf{\Gamma} \in \mathbb{R}^{p \times q}$ measure the time-invariant effects. The random errors e_{it} and $\boldsymbol{\varepsilon}_{it} = (\varepsilon_{it}^1, \dots, \varepsilon_{it}^p)^T$ are zero-mean, where e_{it} is independent of X_{it} , \mathbf{M}_{it} , and \mathbf{Z}_i , and $\boldsymbol{\varepsilon}_{it}$ is independent of X_{it} , $Y_{i,t-1}$, and \mathbf{Z}_i . We define $\boldsymbol{\Theta} = (\boldsymbol{\theta}^T, \boldsymbol{\vartheta}^{1T}, \dots, \boldsymbol{\vartheta}^{pT})^T$ with $\boldsymbol{\theta} = (\alpha_1, \dots, \alpha_T, \beta_1^T, \dots, \beta_T^T, \delta^T)^T$ and $\boldsymbol{\vartheta}^j = (\gamma_1^j, \dots, \gamma_T^j, \eta_2^j, \dots, \eta_T^j, \mathbf{\Gamma}^{jT})^T$, where $\mathbf{\Gamma}^j$ is the j -th row of $\mathbf{\Gamma}$ for $j = 1, \dots, p$.

Figure 2 illustrates the relationships between variables in models (1) and (2) without presenting \mathbf{Z}_i , where solid unidirectional arrows represent causal relationships, while dashed bidirectional arrows indicate correlations. For a given time point t , α_t represents the direct effect from the exposure X_{it} to the outcome Y_{it} , and $\gamma_t^j \beta_t^j$ represents the indirect effect through the mediator M_{it}^j . There are correlations among the mediators $(M_{it}^1, \dots, M_{it}^p)$. In the context of our real data, the DNAm levels of multiple CpG sites exhibit correlations, reflecting their shared regulatory roles or proximity to similar biological functions (Mou et al. 2022).

For adjacent time points $t-1$ and t , we assume both correlation and causation exist among the exposure, mediators, and outcome. We postulate correlations among repeated measurements of outcome Y_{it} . And similarly, we assume correlations among repeated measurements of mediator M_{it}^j . Regarding the causal relationships across time points, since X_{it} reflects the accumulated level of trauma exposure, it is influenced mainly by the exposure at the previous time point. Based on existing literature, DNAm is primarily influenced by external adverse environmental exposures or disease conditions (Turecki and Meaney 2016), we specify the causal effect of PTSD condition $Y_{i,t-1}$ on $(M_{it}^1, \dots, M_{it}^p)$ to be $(\eta_t^1, \dots, \eta_t^p)$. Following the Markov assumption (Luo et al. 2023), we assume that spillover effects among variables occur exclusively at adjacent time points.

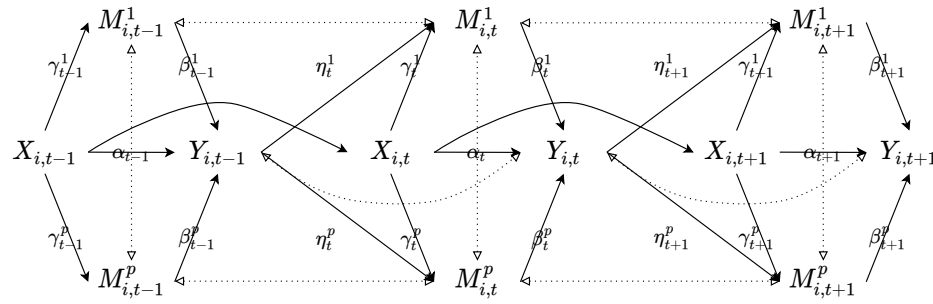


Figure 2: Dynamic mediation structure among time points $t - 1$, t , and $t + 1$. Solid unidirectional arrows represent causal relationships, while dashed bidirectional arrows indicate correlations.

Regression coefficients in models (1) and (2) can be linked to causal estimands. Let $Y^{x,m}$ be the potential outcome that would have been observed if X was set to x and M was set to m , and M^x be the potential mediators that would have been observed if X was set to x . Following Bind et al. (2016), we formalize the following assumptions:

- (i) $Y_{it}^{x,m} \perp\!\!\!\perp X_{it} | \mathbf{Z}_i = \mathbf{z}$, (ii) $M_{it}^x \perp\!\!\!\perp X_{it} | \mathbf{Z}_i = \mathbf{z}$,
- (iii) $Y_{it}^{x,m} \perp\!\!\!\perp M_{it} | X_{it} = x, \mathbf{Z}_i = \mathbf{z}$, (iv) $Y_{it}^{x,m} \perp\!\!\!\perp M_{it}^{\tilde{x}} | \mathbf{Z}_i = \mathbf{z}$,

which indicate that there is no (i) unmeasured exposure-outcome confounding, (ii) exposure-mediator confounding, (iii) mediator-outcome confounding, or (iv) exposure-induced confounding, where \tilde{x} is the realization of exposure at a different value from x . Corresponding to our DNHS data, to estimate the direct effect of trauma exposure X_{it} on PTSD outcome Y_{it} , as well as the indirect effect mediated through DNAm M_{it} , it is essential to control for confounding bias. Assumption (i) requires that all causes of trauma and PTSD are observed and adjusted for. These causes may include demographic information such as age and gender (Wani et al. 2021), which have been accounted for in \mathbf{Z}_i . Assumption (ii) requires that all causes of trauma and DNAm are observed and adjusted for. These may include demographic information like age and gender (Wani et al. 2021), which have been accounted for in \mathbf{Z}_i . Assumption (iii) requires that all causes of DNAm and PTSD are observed and adjusted for. These may include demographic information such as age and gender, as well as blood work (Ocean et al. 2022), which have been accounted for in \mathbf{Z}_i . Assumption (iv) requires that the causes (except for trauma) of DNAm and PTSD are not induced by trauma. This assumption is reasonable, as the demographic information and blood work we control for are intrinsic attributes of individuals. Note that although

$Y_{i,t-1}$ may have a spillover effect on $(M_{it}^1, \dots, M_{it}^p)$, it does not directly affect X_{it} and Y_{it} . Consequently, there is no time-varying confounding in our study, similar to the arguments by VanderWeele and Tchetgen Tchetgen (2017), and thus following Bind et al. (2016), the natural direct effect at time point t for x versus \tilde{x} is $\mathbb{E}(Y_{it}^{x, \mathbf{m}_{it}^{\tilde{x}}} - Y_{it}^{\tilde{x}, \mathbf{m}_{it}^{\tilde{x}}}) = \alpha_t(x - \tilde{x})$, and the natural indirect effect at time point t for x versus \tilde{x} is $\mathbb{E}(Y_{it}^{x, \mathbf{m}_{it}^x} - Y_{it}^{x, \mathbf{m}_{it}^{\tilde{x}}}) = \beta_t^T \gamma_t(x - \tilde{x})$.

3.2 Non-monotone missing mediators

In our DNHS data, exposure X_{it} , outcome Y_{it} , and confounders \mathbf{Z}_i are always observed, but the mediators \mathbf{M}_{it} suffer from non-monotone missingness due to the following reasons: $\mathbf{M}_i = (\mathbf{M}_{i1}, \dots, \mathbf{M}_{iT})^T$ of some subjects may be available only at certain visits, but missing at the next time point, and measured again at later visits. Given the non-monotone missingness, subjects can be divided into R disjoint groups based on different missingness patterns across all time points. Figure 3 illustrates an example where we have $T = 3$ time points and $R = 7$ groups. Each blank column represents unobserved mediators, while the colored ones represent observed mediators. Notice that subjects in group 1 have observed \mathbf{M}_i at all time points, but subjects in group 2 have observed \mathbf{M}_i only at time points 1 and 2.

Define the missing indicator $\mathbf{R}_{it} = (R_{it}^1, \dots, R_{it}^p)^T$, where $R_{it}^j = 1$ for $j = 1, \dots, p$ if \mathbf{M}_{it} is observed at time point t , and $R_{it}^j = 0$ for $j = 1, \dots, p$ if \mathbf{M}_{it} is missing at time point t . In this paper, we assume that the missingness mechanism is missing at random (MAR, Little and Rubin 2019), meaning that the conditional distribution of the missing indicator $\mathbf{R}_i = (\mathbf{R}_{i1}, \dots, \mathbf{R}_{iT})^T$, denoted as $f_{\mathbf{R}_i}(\mathbf{r}_i | \mathbf{X}_i, \mathbf{Y}_i, \mathbf{Z}_i, \mathbf{M}_i)$, depends only on the exposure $\mathbf{X}_i = (X_{i1}, \dots, X_{iT})^T$, the outcome $\mathbf{Y}_i = (Y_{i1}, \dots, Y_{iT})^T$, the confounders \mathbf{Z}_i , and the observed mediator component $\mathbf{R}_i \mathbf{M}_i$ (Wei et al. 2022).

Let $\mathcal{H}(r)$ be the index set of subjects in group r and $n_r = |\mathcal{H}(r)|$ be the sample size ($r = 1, \dots, R$). Let $o(r)$ and $u(r)$ be the index sets of time points corresponding to observed and unobserved mediators in group r . Consequently, $\mathbf{M}_{o(r)}$ and $\mathbf{M}_{u(r)}$ represent observed and unobserved mediators in group r . In addition, for each $t \in u(r)$, we define $\mathcal{T}(r, t)$ as the collection of time sets, and each element in $\mathcal{T}(r, t)$ is a subset of $o(r)$. The $\mathcal{T}(r, t)$ can be obtained by the following steps: (i) let $\mathcal{G}(r, t)$ be the index set of the groups where mediators are observed not only at time point t , but also at least one time point in $o(r)$; (ii) for each $r' \in \mathcal{G}(r, t)$, let $\mathcal{T}_{r'} = o(r) \cap o(r')$; (iii) $\mathcal{T}(r, t) = \bigcup_{r' \in \mathcal{G}(r, t)} \mathcal{T}_{r'}$. Note that the

definition of $\mathcal{T}(r, t)$ is to prepare for the imputation method later. That is, we can use observed values at each time set $t \in \mathcal{T}(r, t)$ to impute the unobserved values at time point $t \in u(r)$.

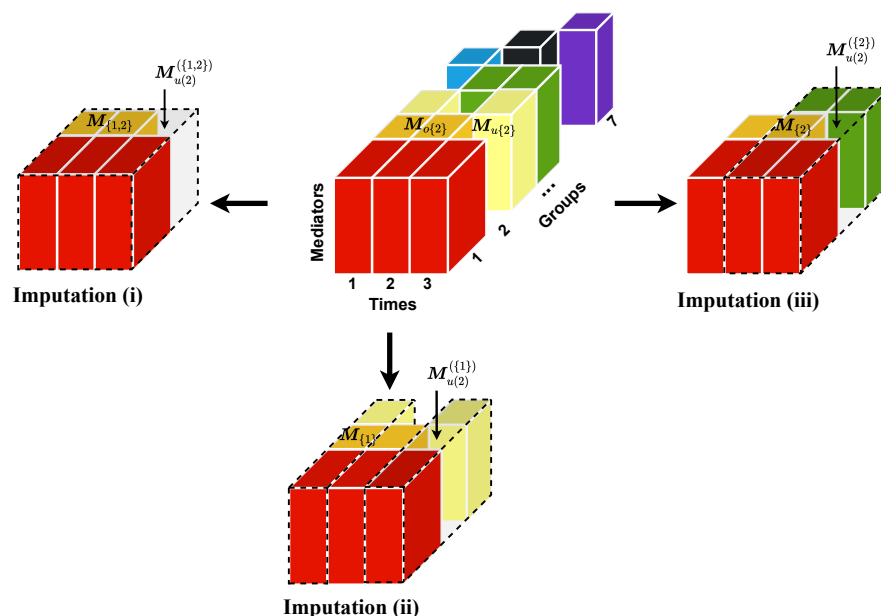


Figure 3: Middle: Non-monotone missingness structure of mediators across all time points. Each blank column represents unobserved mediators, while the colored ones represent observed mediators. Surroundings: Multiple imputations for the missing column $M_{u(2)}$ in group 2.

4 Methods

In this section, we first propose a new imputation method called Longitudinal Multiple Imputation (LMI) to handle non-monotone missing data. Second, we adopt the generalized method of moments (GMM, Hansen 1982) to combine the multiple imputations. Third, we propose a regularization approach to identify relevant mediators and estimate time-varying effects. Finally, we employ principal component analysis to solve the singularity issue of the weighting matrix in GMM.

4.1 Longitudinal multiple imputation (LMI)

In this subsection, we propose the LMI approach to handle non-monotone missing longitudinal data. The main idea is that, due to the correlation between repeated measurements \mathbf{M}_t ($t = 1, \dots, T$) from the same subject, we can use the observed values at some time points to predict the unobserved values at other time points. For example, for subjects with observed \mathbf{M}_t at time points 1 and 2, we can use some statistical methods, such as regression, to learn the correlation of variables between these two time points. Thus, for some other individuals with observed \mathbf{M}_t at time point 1 but missing \mathbf{M}_t at time point 2, we can impute the data at time point 2 based on the learned correlation.

Given a group r with a time point $t \in u(r)$, subjects have unobserved M_t^j at this time point, but have observed \mathbf{M}_t at each time set $\mathcal{t} \in \mathcal{T}(r, t)$. We can use observed values at each time set $\mathcal{t} \in \mathcal{T}(r, t)$ to impute the unobserved values at time point $t \in u(r)$. Specifically, given a time set $\mathcal{t} \in \mathcal{T}(r, t)$, we first estimate $\mathbb{E}(M_t^j | \mathbf{M}_t, \mathbf{D})$ using other groups containing observed M_t^j and \mathbf{M}_t , where \mathbf{D} contains exposure $\mathbf{X} = (X_1, \dots, X_T)^T$ at all time points, outcome Y_{t-1} at time point $t - 1$, and confounders \mathbf{Z} . Based on the learned correlation, we then impute unobserved M_t^j using \mathbf{M}_t and \mathbf{D} . Let $\mathbf{M}_t^{(\mathcal{t})} = \{\mathbb{E}(M_t^j | \mathbf{M}_t, \mathbf{D}), j = 1, \dots, p\}$ represent the imputed values based on the time set $\mathcal{t} \in \mathcal{T}(r, t)$, then we can obtain $|\mathcal{T}(r, t)|$ imputations.

We illustrate the LMI with an example in Figure 3. Specifically, $\mathcal{R}(2)$ refers to the index set of subjects in group 2. Since subjects in group 2 have observed mediators at time points 1 and 2, but have unobserved mediators at time point 3, $o(2) = \{1, 2\}$ and $\mathbf{M}_{o(2)}$ refers to observed mediators, $u(2) = \{3\}$ and $\mathbf{M}_{u(2)}$ refers to unobserved mediators. The collection of time sets $\mathcal{T}(2, 3) = \{\{1, 2\}, \{1\}, \{2\}\}$ since: (i) $\mathcal{G}(2, 3) = \{1, 3, 4\}$; (ii) $\mathcal{t}_1 = o(2) \cap o(1) = \{1, 2\}$, $\mathcal{t}_3 = o(2) \cap o(3) = \{1\}$, and $\mathcal{t}_4 = o(2) \cap o(4) = \{2\}$; (iii) $\mathcal{T}(2, 3) = \bigcup_{r' \in \mathcal{G}(2, 3)} \mathcal{t}_{r'} = \{\{1, 2\}, \{1\}, \{2\}\}$.

Given group 2 with unobserved \mathbf{M}_t at time point $t \in u(r) = 3$, we can impute the unobserved values in $|\mathcal{T}(r, t)| = |\{\{1, 2\}, \{1\}, \{2\}\}| = 3$ ways: for $j = 1, \dots, p$, (i) for $\mathcal{t} = \{1, 2\}$, we estimate $\mathbb{E}(M_3^j | \mathbf{M}_{\{1, 2\}}, \mathbf{D})$ using group 1, and then obtain the imputed values $\mathbf{M}_3^{\{1, 2\}}$; (ii) for $\mathcal{t} = \{1\}$, we estimate $\mathbb{E}(M_3^j | \mathbf{M}_{\{1\}}, \mathbf{D})$ using groups 1 and 3, and then obtain the imputed values $\mathbf{M}_3^{\{1\}}$; (iii) for $\mathcal{t} = \{2\}$, we estimate $\mathbb{E}(M_3^j | \mathbf{M}_{\{2\}}, \mathbf{D})$ using groups 1 and 4, and then obtain the imputed values $\mathbf{M}_3^{\{2\}}$.

Remark 1. One novelty of our LMI method is to make full use of the dependence among

repeated measurements, which exploits the inherent correlation structure of longitudinal data, not based on stringent low-rank assumptions (Mazumder, Hastie and Tibshirani 2010) or information-poor trajectory means (Jahangiri et al. 2023). The LMI method allows high-dimensional variables to be entirely unobserved at some time points, while existing approaches only consider missingness of the outcome or with a few covariates (Chen et al. 2021). Furthermore, the LMI method allows for the absence of individuals with complete repeated measurements in the dataset, whereas existing approaches may require that there must be some subjects with complete data, or that the data at the first time visit is always observed for all subjects (Jahangiri et al. 2023).

4.2 Combining multiple imputations

In this subsection, we propose to combine the information from all LMIs. Since different imputations correspond to different estimating functions, based on moment conditions (Hansen 1982), we integrate all information through the GMM.

First, we construct the moment conditions corresponding to model (1). For a given group r , $j = 1, \dots, p$, and each time set $\mathcal{t} \in \mathcal{T}(r, t)$, since unobserved M_t^j at time point $t \in u(r)$ is imputed through $M_t^{j(\mathcal{t})} = \mathbb{E}(M_t^j | \mathbf{M}_t, \mathbf{D})$, the covariates \mathbf{M}_t and \mathbf{D} are uncorrelated with residuals of the projection $M_t^j - M_t^{j(\mathcal{t})}$. Therefore,

$$\begin{aligned} & \mathbb{E} \left\{ \left(X_t, \mathbf{M}_t^{(\mathcal{t})T} \right)^T \left(Y_t - \alpha_{0t} X_t - \beta_{0t}^T \mathbf{M}_t^{(\mathcal{t})} - \delta_0^T \mathbf{Z} \right) \right\} \\ &= \mathbb{E} \left\{ \left(X_t, \mathbf{M}_t^{(\mathcal{t})T} \right)^T \left(e_t + \beta_{0t}^T \left(\mathbf{M}_t - \mathbf{M}_t^{(\mathcal{t})} \right) \right) \right\} = \mathbf{0}, \end{aligned}$$

and

$$\begin{aligned} & \mathbb{E} \left\{ \sum_{t=1}^T \mathbf{Z} \left(Y_t - \alpha_{0t} X_t - \beta_{0t}^T \mathbf{M}_t^{(\mathcal{t})} - \delta_0^T \mathbf{Z} \right) \right\} \\ &= \mathbb{E} \left\{ \sum_{t=1}^T \mathbf{Z} \left(e_t + \beta_{0t}^T \left(\mathbf{M}_t - \mathbf{M}_t^{(\mathcal{t})} \right) \right) \right\} = \mathbf{0}. \end{aligned}$$

The moment conditions are satisfied for the true parameter $\boldsymbol{\theta}_0 = (\alpha_{01}, \dots, \alpha_{0T}, \beta_{01}^T, \dots, \beta_{0T}^T, \delta_0^T)^T$.

We then construct estimating functions based on these moment conditions. For each $i \in \mathcal{X}(r)$, let $\mathbf{M}_{it}^{(\mathcal{t})}$ be the imputed values based on time set $\mathcal{t} \in \mathcal{T}(r, t)$, where $\mathbf{M}_{it}^{(\mathcal{t})} = \mathbf{M}_{it}$ for $t \in o(r)$, and $M_{it}^{j(\mathcal{t})} = \mathbb{E}(M_{it}^j | \mathbf{M}_{it}, \mathbf{D}_i)$ for $t \in u(r)$ and $j = 1, \dots, p$. The estimating

functions for X_{it} and $\mathbf{M}_{it}^{(\ell)}$ at time point $t \in o(r)$ are $\mathbf{g}_{it}^{(r)}(\boldsymbol{\theta}) = (X_{it}, \mathbf{M}_{it}^{(\ell)T})^T (Y_{it} - \alpha_t X_{it} - \beta_t^T \mathbf{M}_{it}^{(\ell)} - \delta^T \mathbf{Z}_i)$ in which $\mathbf{M}_{it}^{(\ell)} = \mathbf{M}_{it}$, and the estimating functions for X_{it} and $\mathbf{M}_{it}^{(\ell)}$ at time point $t \in u(r)$ are $\mathbf{g}_{it}^{(r)}(\boldsymbol{\theta}) = \{\mathbf{g}_{it}^{(r,\ell)}(\boldsymbol{\theta}), \ell \in \mathcal{T}(r, t)\}$, where

$$\mathbf{g}_{it}^{(r,\ell)}(\boldsymbol{\theta}) = \left(X_{it}, \mathbf{M}_{it}^{(\ell)T} \right)^T \left(Y_{it} - \alpha_t X_{it} - \beta_t^T \mathbf{M}_{it}^{(\ell)} - \delta^T \mathbf{Z}_i \right). \quad (3)$$

The estimating functions for \mathbf{Z}_i are $\mathbf{g}_i^{(r)}(\boldsymbol{\theta}) = \{\sum_{t=1}^T \mathbf{g}_{it}^{(r)}(\boldsymbol{\theta})\}$, where $\mathbf{g}_{it}^{(r)}(\boldsymbol{\theta}) = \mathbf{Z}_i (Y_{it} - \alpha_t X_{it} - \beta_t^T \mathbf{M}_{it}^{(\ell)} - \delta^T \mathbf{Z}_i)$ at time point $t \in o(r)$ in which $\mathbf{M}_{it}^{(\ell)} = \mathbf{M}_{it}$, and $\mathbf{g}_{it}^{(r)}(\boldsymbol{\theta}) \in \{\mathbf{g}_{it}^{(r,\ell)}(\boldsymbol{\theta}), \ell \in \mathcal{T}(r, t)\}$ at time point $t \in u(r)$ in which

$$\mathbf{g}_{it}^{(r,\ell)}(\boldsymbol{\theta}) = \mathbf{Z}_i \left(Y_{it} - \alpha_t X_{it} - \beta_t^T \mathbf{M}_{it}^{(\ell)} - \delta^T \mathbf{Z}_i \right). \quad (4)$$

To integrate information from all groups, we propose an aggregated vector of estimating functions:

$$\mathbf{g}(\boldsymbol{\theta}) = \left\{ (\mathbf{g}^{(1)}(\boldsymbol{\theta}))^T, \dots, (\mathbf{g}^{(R)}(\boldsymbol{\theta}))^T \right\}^T, \quad (5)$$

where $\mathbf{g}^{(r)}(\boldsymbol{\theta}) = \frac{1}{n_r} \sum_{i \in \mathcal{H}(r)} \{(\mathbf{g}_{i1}^{(r)}(\boldsymbol{\theta}))^T, \dots, (\mathbf{g}_{iT}^{(r)}(\boldsymbol{\theta}))^T, (\mathbf{g}_i^{(r)}(\boldsymbol{\theta}))^T\}^T$. Note that the total number of estimating functions exceeds the dimension of parameters, and the estimating functions from groups with fewer missingness or more precise imputations tend to have smaller variance. To combine all estimating functions in $\mathbf{g}(\boldsymbol{\theta})$ and put more weights on high-quality moments, we propose the following quadratic loss function:

$$G(\boldsymbol{\theta}) = \sum_{r=1}^R (\mathbf{g}^{(r)}(\boldsymbol{\theta}))^T (\mathbf{W}^{(r)}(\boldsymbol{\theta}))^{-1} \mathbf{g}^{(r)}(\boldsymbol{\theta}), \quad (6)$$

with weighting matrix $\mathbf{W}^{(r)}(\boldsymbol{\theta}) = \frac{1}{n_r} \sum_{i \in \mathcal{H}(r)} \mathbf{g}^{(r)}(\boldsymbol{\theta}) (\mathbf{g}^{(r)}(\boldsymbol{\theta}))^T$.

For each mediator M^j ($j = 1, \dots, p$) in model (2), we can get the aggregated vector of estimating functions $\mathbf{h}^j(\boldsymbol{\vartheta}^j)$ similar to (5). Detailed formulations are shown in the Supplementary Material. We propose the following quadratic loss function:

$$H^j(\boldsymbol{\vartheta}^j) = \sum_{r=1}^R (\mathbf{h}^{j(r)}(\boldsymbol{\vartheta}^j))^T (\mathbf{W}^{j(r)}(\boldsymbol{\vartheta}^j))^{-1} \mathbf{h}^{j(r)}(\boldsymbol{\vartheta}^j), \quad (7)$$

with weighting matrix $\mathbf{W}^{j(r)}(\boldsymbol{\vartheta}^j) = \frac{1}{n_r} \sum_{i \in \mathcal{H}(r)} \mathbf{h}^{j(r)}(\boldsymbol{\vartheta}^j) (\mathbf{h}^{j(r)}(\boldsymbol{\vartheta}^j))^T$. The loss function incorporating models (1) and (2) for all subjects is $G(\boldsymbol{\theta}) + \sum_{j=1}^p H^j(\boldsymbol{\vartheta}^j)$.

Remark 2. For multiple imputed data, different imputations correspond to different estimating functions, all of which satisfy the moment conditions. We use GMM to combine all estimating functions, which adaptively puts more weights to imputations with smaller

variance, while still consider the information contained in imputations with larger variance. In contrast, if we only use the imputation with the smallest variance while discarding the others, it may lead to efficiency losses due to not use of available information.

4.3 Regularization

To select the relevant mediators from high-dimensional variables and to estimate the time-varying effects, we adopt the regularization approach. Specifically, we propose to incorporate three penalty terms to: (i) identify mediators with large mediation effects at each time point, (ii) shrink similar effects at adjacent time points, and (iii) shrink effects of demographic variables. Consequently, the objective function of the proposed method is

$$L(\Theta) = G(\theta) + \sum_{j=1}^p H^j(\vartheta^j) + P_1(\Theta) + P_2(\Theta) + P_3(\Theta), \quad (8)$$

where

$$P_1(\Theta) = \sum_{t=1}^T P_{\text{SCAD}, \lambda_1, a}(|\alpha_t|) + \lambda_1 \sum_{t=1}^T \sum_{j=1}^p \left\{ 1 - \frac{1}{(1+b|\beta_t^j|)(1+b|\gamma_t^j|)} \right\}, \quad (9)$$

is a penalty for direct effects and indirect effects, with $P'_{\text{SCAD}, \lambda_1, a}(\alpha) = \lambda_1 \{\mathbb{I}(\alpha \leq \lambda_1) + \frac{(a\lambda_1 - \alpha)_+}{(a-1)\lambda_1} \mathbb{I}(\alpha > \lambda_1)\}$ for some $\alpha > 0$, $a > 2$, and $b > 0$,

$$P_2(\Theta) = \lambda_2 \sum_{t=2}^T |\alpha_t - \alpha_{t-1}| + \lambda_2 \sum_{t=2}^T \sum_{j=1}^p (|\beta_t^j - \beta_{t-1}^j| + |\gamma_t^j - \gamma_{t-1}^j|), \quad (10)$$

is a fusion penalty for time-varying effects, and

$$P_3(\Theta) = \lambda_3 \left(\|\delta\|_1 + \sum_{j=1}^p \|\Gamma^j\|_1 \right), \quad (11)$$

is a penalty for demographic effects. The regularization parameters λ_1 , λ_2 , and λ_3 control the amount of shrinkage.

The second term in $P_1(\Theta)$ is a bi-level mediation penalty for indirect effects, which links models (1) and (2) by penalizing β_t and γ_t jointly rather than separately, and achieves bi-level variable selection (Huang et al. 2009). This estimator not only obtains the sparse solution for β_t^j and γ_t^j when $\beta_t^j = 0$ and $\gamma_t^j = 0$, but also obtains the sparse solution for β_t^j or γ_t^j when $\beta_t^j = 0$ or $\gamma_t^j = 0$. On the other hand, it inherits the advantage of the smoothly clipped absolute deviation (SCAD, Fan and Li 2001) penalty where the function gradually

levels off as β_t^j and γ_t^j increase, which avoids over-shrinkage. In general, it is possible that not every mediator has different effects at different time points. We use the fused lasso penalty (Tibshirani et al. 2005) in $P_2(\Theta)$ to shrink similar effects at adjacent time points, which links β_t and γ_t across all time points. $P_3(\Theta)$ is a lasso penalty (Friedman, Hastie and Tibshirani 2010) for selecting demographic effects to avoid over-fitting.

4.4 Principal components of weighting matrix

The GMM estimators in (6) and (7) may have poor finite sample performance in highly overidentified models, due to imprecise estimation of the weighting matrix $\mathbf{W}^{(r)}(\theta)$ or $\mathbf{W}^{j(r)}(\vartheta^j)$, or under the singularity of the weighting matrix (Doran and Schmidt 2006). This could be attributed to the large number of estimating functions compared to the relatively small sample size, or due to the overlap of information between LMIs. For example, as illustrated in Figure 3, the observed values at time point 1 in group 2 are used for the estimation of both $\mathbf{M}_{u(2)}^{\{1,2\}}$ and $\mathbf{M}_{u(2)}^{\{1\}}$, and the observed values at time point 2 in group 2 are used for the estimation of both $\mathbf{M}_{u(2)}^{\{1,2\}}$ and $\mathbf{M}_{u(2)}^{\{2\}}$.

To address the singularity issue, we reduce the dimension of $\mathbf{g}^{(r)}$ and $\mathbf{h}^{j(r)}$ ($r = 1, \dots, R$ and $j = 1, \dots, p$), by combining informative estimating functions, for example, using the first several largest principal components. If $\mathbf{W}^{(r)} = \frac{1}{n_r} \sum_{i \in \mathcal{G}(r)} \mathbf{g}_i^{(r)} \mathbf{g}_i^{(r)T}$ is singular or close to singular, we extract the first $u^{(r)}$ principal components $\tilde{\mathbf{g}}^{(r)} = \mathbf{U}^{(r)} \mathbf{g}^{(r)}$ from $\mathbf{g}^{(r)}$, where $\mathbf{U}^{(r)}$ contains $u^{(r)}$ eigenvectors of $\mathbf{W}^{(r)}$ corresponding to the largest $u^{(r)}$ nonzero eigenvalues. Similarly, if $\mathbf{W}^{j(r)} = \frac{1}{n_r} \sum_{i \in \mathcal{G}(r)} \mathbf{h}_i^{j(r)} \mathbf{h}_i^{j(r)T}$ is singular or close to singular, we extract the first $u^{j(r)}$ principal components $\tilde{\mathbf{h}}^{j(r)} = \mathbf{U}^{j(r)} \mathbf{h}^{j(r)}$ from $\mathbf{h}^{j(r)}$, where $\mathbf{U}^{j(r)}$ contains $u^{j(r)}$ eigenvectors of $\mathbf{W}^{j(r)}$ corresponding to the largest $u^{j(r)}$ nonzero eigenvalues. The numbers of principal components $u^{(r)}$ and $u^{j(r)}$ ($r = 1, \dots, R$ and $j = 1, \dots, p$) can be selected based on the information criterion proposed by Cho and Qu (2015) to capture sufficient information from the estimating functions. Note that if $\mathbf{W}^{(r)}$ or $\mathbf{W}^{j(r)}$ is not singular, $\mathbf{U}^{(r)}$ or $\mathbf{U}^{j(r)}$ is an identity matrix. Consequently, the final objective function is

$$\tilde{L}(\Theta) = \tilde{G}(\theta) + \sum_{j=1}^p \tilde{H}^j(\vartheta^j) + P_1(\Theta) + P_2(\Theta) + P_3(\Theta), \quad (12)$$

where the quadratic form $\tilde{G}(\theta) = \sum_{r=1}^R (\tilde{\mathbf{g}}^{(r)}(\theta))^T (\tilde{\mathbf{W}}^{(r)}(\theta))^{-1} \tilde{\mathbf{g}}^{(r)}(\theta)$ with transformed weighting matrix $\tilde{\mathbf{W}}^{(r)}(\theta) = \frac{1}{n_r} \sum_{i \in \mathcal{R}(r)} \tilde{\mathbf{g}}^{(r)}(\theta) (\tilde{\mathbf{g}}^{(r)}(\theta))^T$, and the quadratic form $\tilde{H}^j(\vartheta^j) =$

$$\sum_{r=1}^R (\tilde{\mathbf{h}}^{j(r)}(\boldsymbol{\vartheta}^j))^T (\tilde{\mathbf{W}}^{j(r)}(\boldsymbol{\vartheta}^j))^{-1} \tilde{\mathbf{h}}^{j(r)}(\boldsymbol{\vartheta}^j) \text{ with transformed weighting matrix } \tilde{\mathbf{W}}^{j(r)}(\boldsymbol{\vartheta}^j) = \frac{1}{n_r} \sum_{i \in \mathcal{H}(r)} \tilde{\mathbf{h}}^{j(r)}(\boldsymbol{\vartheta}^j) (\tilde{\mathbf{h}}^{j(r)}(\boldsymbol{\vartheta}^j))^T.$$

5 Computation

In this section, we provide the computational details of the proposed method. We first describe the implementation of LMI from Section 4.1. Then we introduce an information criterion to select the numbers of principal components from Section 4.4. Next, we present a smoothing technique to handle the non-separable fusion penalty from Section 4.3. Afterwards, we introduce the overall optimization algorithm of the proposed method for a given set of regularization parameters $(\lambda_1, \lambda_2, \lambda_3)^T$ from Section 4.3. Finally, we propose an information criterion to select these regularization parameters.

When conducting LMI, the conditional expectation $\mathbb{E}(M_t^j | \mathbf{M}_t, \mathbf{D})$ is required to be estimated for each time set $t \in \mathcal{T}(r, t)$, time point $t \in u(r)$, group $r = 1, \dots, R$, and dimension $j = 1, \dots, p$. In this paper, we adopt the lasso-regularized linear model (Friedman, Hastie and Tibshirani 2010) to account for high-dimensionality, where the regularization parameter is selected by cross-validation.

To capture most of the information from the estimating functions $\mathbf{g}^{(r)}$ and $\mathbf{h}^{j(r)}$ ($r = 1, \dots, R$ and $j = 1, \dots, p$), we select the numbers of principal components $u^{(r)}$ and $u^{j(r)}$, by minimizing the following Bayesian-type information criterion (Cho and Qu 2015):

$$\Psi(u) = \frac{\text{tr}\{\boldsymbol{\Omega} - \tilde{\boldsymbol{\Omega}}(u)\}}{\text{tr}\{\boldsymbol{\Omega}\}} + u \frac{\log(n_r d)}{n_r d}, \quad (13)$$

where $\text{tr}\{\cdot\}$ is the trace of a matrix, $\boldsymbol{\Omega}$ represents the weighting matrix $\mathbf{W}^{(r)}$ or $\mathbf{W}^{j(r)}$ with dimension d , $\tilde{\boldsymbol{\Omega}}(u) = \sum_{k=1}^u \xi_k \boldsymbol{\nu}_k \boldsymbol{\nu}_k^T$ is an approximation of $\boldsymbol{\Omega}$ based on the largest u eigenvectors, and ξ_k is the k -th largest eigenvalue of $\boldsymbol{\Omega}$ corresponding to the eigenvector $\boldsymbol{\nu}_k$. Since $\text{tr}\{\boldsymbol{\Omega} - \tilde{\boldsymbol{\Omega}}(u)\} = \sum_{k=u+1}^d \xi_k$, the minimizer of $\Psi(u)$ is indeed the number of eigenvalues greater than $\text{tr}\{\boldsymbol{\Omega}\} \log(n_r d) / (n_r d)$.

Since the non-smooth and non-separable of fusion penalty $P_2(\boldsymbol{\Theta})$ in (10) makes the optimization challenging, we use the smoothing technique from Chen et al. (2012). Specifically, we reformulate the fused lasso as $P_2(\boldsymbol{\Theta}) = \lambda_2 \|\mathbb{D}\boldsymbol{\Theta}\|_1 = \lambda_2 \max_{\|\boldsymbol{\eta}\|_\infty \leq 1} (\boldsymbol{\eta}^T \mathbb{D}\boldsymbol{\Theta})$, where \mathbb{D} is a difference operator corresponding to the differences in $P_2(\boldsymbol{\Theta})$. Let $\tilde{P}_2(\boldsymbol{\Theta}; \rho) =$

$\lambda_2 \max_{\|\boldsymbol{\eta}\|_\infty \leq 1} (\boldsymbol{\eta}^T \mathbb{D} \boldsymbol{\Theta} - \rho \|\boldsymbol{\eta}\|_2^2 / 2)$, where ρ is a positive smoothing parameter and $\tilde{P}_2(\boldsymbol{\Theta}; \rho)$ approximates $P_2(\boldsymbol{\Theta})$ as $\rho \rightarrow 0$. Define a projection operator $S(x) = -\mathbb{I}(x < -1) + x\mathbb{I}(-1 \leq x \leq 1) + \mathbb{I}(x > 1)$, we can deduce that $\boldsymbol{\eta}^* = S(\mathbb{D} \boldsymbol{\Theta} / \rho)$ is indeed the optimal solution in $\tilde{P}_2(\boldsymbol{\Theta}; \rho)$. Thus we have

$$\tilde{P}_2(\boldsymbol{\Theta}, \rho) = \lambda_2 \left(\boldsymbol{\eta}^{*T} \mathbb{D} \boldsymbol{\Theta} - \frac{\rho}{2} \|\boldsymbol{\eta}^*\|_2^2 \right), \quad (14)$$

which is convex and differentiable in $\boldsymbol{\Theta}$. We set $\rho = 10^{-4}$ following [Chen et al. \(2012\)](#).

We propose an algorithm based on the iterated GMM with principal component analysis. Specifically, given $\boldsymbol{\Theta}^{(s-1)}$ from the $(s-1)$ -th iteration, for $r = 1, \dots, R$ and $j = 1, \dots, p$, if the weighting matrix $\mathbf{W}^{(r)}(\boldsymbol{\theta}^{(s-1)})$ or $\mathbf{W}^{j(r)}(\boldsymbol{\vartheta}^{j(s-1)})$ is singular or close to singular, we first select the number of principal components $u^{(r)(s-1)}$ and $u^{j(r)(s-1)}$ by minimizing the information criterion in (13), then construct principal component matrices $\mathbf{U}^{(r)}(\boldsymbol{\theta}^{(s-1)})$ and $\mathbf{U}^{j(r)}(\boldsymbol{\vartheta}^{j(s-1)})$, next get the reduced moment conditions $\tilde{\mathbf{g}}^{(r)}(\boldsymbol{\theta}^{(s-1)})$ and $\tilde{\mathbf{h}}^{j(r)}(\boldsymbol{\vartheta}^{j(s-1)})$, and finally obtain the transformed weighting matrices $\tilde{\mathbf{W}}^{(r)}(\boldsymbol{\theta}^{(s-1)})$ and $\tilde{\mathbf{W}}^{j(r)}(\boldsymbol{\vartheta}^{j(s-1)})$.

At the s -th iteration, we solve

$$\boldsymbol{\Theta}^{(s)} = \underset{\boldsymbol{\Theta}}{\operatorname{argmin}} \tilde{L}^{(s)}(\boldsymbol{\Theta}), \quad (15)$$

where $\tilde{L}^{(s)}(\boldsymbol{\Theta}) = \tilde{G}^{(s-1)}(\boldsymbol{\theta}) + \sum_{j=1}^p \tilde{H}^{j(s-1)}(\boldsymbol{\vartheta}^j) + P_1(\boldsymbol{\Theta}) + \tilde{P}_2(\boldsymbol{\Theta}; \rho) + P_3(\boldsymbol{\Theta})$. The quadratic form $\tilde{G}^{(s-1)}(\boldsymbol{\theta}) = \sum_{r=1}^R (\tilde{\mathbf{g}}^{(r)(s-1)}(\boldsymbol{\theta}))^T (\tilde{\mathbf{W}}^{(r)}(\boldsymbol{\theta}^{(s-1)}))^{-1} \tilde{\mathbf{g}}^{(r)(s-1)}(\boldsymbol{\theta})$ with $\tilde{\mathbf{g}}^{(r)(s-1)}(\boldsymbol{\theta}) = \mathbf{U}^{(r)}(\boldsymbol{\theta}^{(s-1)}) \mathbf{g}^{(r)}(\boldsymbol{\theta})$ uses principal component matrix $\mathbf{U}^{(r)}(\boldsymbol{\theta}^{(s-1)})$ and transformed weighting matrix $\tilde{\mathbf{W}}^{(r)}(\boldsymbol{\theta}^{(s-1)})$ obtained at the $(s-1)$ -th iteration. The quadratic form $\tilde{H}^{j(s-1)}(\boldsymbol{\vartheta}^j) = \sum_{r=1}^R (\tilde{\mathbf{h}}^{j(r)(s-1)}(\boldsymbol{\vartheta}^j))^T (\tilde{\mathbf{W}}^{j(r)}(\boldsymbol{\vartheta}^{j(s-1)}))^{-1} \tilde{\mathbf{h}}^{j(r)(s-1)}(\boldsymbol{\vartheta}^j)$ uses principal component matrix $\mathbf{U}^{j(r)}(\boldsymbol{\vartheta}^{j(s-1)})$ and transformed weighting matrix $\tilde{\mathbf{W}}^{j(r)}(\boldsymbol{\vartheta}^{j(s-1)})$ obtained at the $(s-1)$ -th iteration. The above calculation process can be summarized in Algorithm 1, and an expanded version of this algorithm is shown in the Supplementary Material.

The turning parameter a in (9) is set to be 3.7 following [Fan and Li \(2001\)](#) and b in (9) is set to be 10 following [Xue et al. \(2022\)](#). To determine the regularization parameters $\boldsymbol{\lambda} = (\lambda_1, \lambda_2, \lambda_3)^T$, we propose the following Bayesian-type information criterion ([Xue and Qu 2021](#), [Xue et al. 2022](#)):

$$\Phi(\lambda_1, \lambda_2, \lambda_3) = \log \left\{ \frac{\text{RSS}_Y(\hat{\boldsymbol{\Theta}}_{\boldsymbol{\lambda}})}{NT} \right\} + \log \left\{ \sum_{j=1}^p \frac{\text{RSS}_{M^j}(\hat{\boldsymbol{\Theta}}_{\boldsymbol{\lambda}})}{NT} \right\} + \frac{\log(NT)}{NT} \text{df}_{\boldsymbol{\lambda}}, \quad (16)$$

Algorithm 1

- 1: LMI. Conduct longitudinal multiple imputation.
 - 2: Initialization. Obtain initial values $\Theta^{(0)}$ based on the averaged imputed data. Set the regularization parameters λ_1 , λ_2 , and λ_3 .
 - 3: At the s -th iteration, given $\Theta^{(s-1)}$ from the $(s-1)$ -th iteration:
 - (a) Principal component analysis. For $r = 1, \dots, R$ and $j = 1, \dots, p$, select the number of principal components using (13) if $\mathbf{W}^{(r)}(\boldsymbol{\theta}^{(s-1)})$ or $\mathbf{W}^{j(r)}(\boldsymbol{\vartheta}^{j(s-1)})$ is singular or close to singular, and obtain $\mathbf{U}^{(r)}(\boldsymbol{\theta}^{(s-1)})$, $\mathbf{U}^{j(r)}(\boldsymbol{\vartheta}^{j(s-1)})$, $\tilde{\mathbf{W}}^{(r)}(\boldsymbol{\theta}^{(s-1)})$, and $\tilde{\mathbf{W}}^{j(r)}(\boldsymbol{\vartheta}^{j(s-1)})$.
 - (b) Optimization. Solve (15) through the gradient descent to obtain $\Theta^{(s)}$.
 - 4: Iterate Step 3 until $\|\Theta^{(s)} - \Theta^{(s-1)}\|_1 < 10^{-6}$. Return $\Theta^{(s)}$.
-

where $\text{RSS}_Y(\hat{\Theta}_\lambda) = \sum_{r=1}^R \sum_{i \in \mathcal{H}(r)} \sum_{t=1}^T \frac{1}{|\mathcal{T}(r,t)|} \sum_{t \in \mathcal{T}(r,t)} (Y_{it} - \hat{\alpha}_t X_{it} - \hat{\beta}_t^T \mathbf{M}_{it}^{(t)} - \hat{\delta}^T \mathbf{Z}_i)^2$ and $\text{RSS}_{M^j}(\hat{\Theta}_\lambda) = \sum_{r=1}^R \sum_{i \in \mathcal{H}(r)} \sum_{t=1}^T \frac{1}{|\mathcal{T}(r,t)|} \sum_{t \in \mathcal{T}(r,t)} (M_{it}^{j(t)} - \hat{\gamma}_t^j X_{it} - \hat{\Gamma}^{jT} \mathbf{Z}_i)^2$ are the residual sum of squares based on imputed data corresponding to models (1) and (2). Degrees of freedom $\text{df}_\lambda = \text{df}_\lambda(\hat{\alpha}) + \text{df}_\lambda(\hat{\beta}) + \text{df}_\lambda(\hat{\gamma}) + \text{df}_\lambda(\hat{\delta}, \hat{\Gamma})$, where $\text{df}_\lambda(\hat{\delta}, \hat{\Gamma})$ is the number of nonzero estimated values in $\hat{\delta}$ and $\hat{\Gamma}$, and

$$\begin{aligned} \text{df}_\lambda(\hat{\alpha}) &= T - \sum_{t=1}^T \mathbb{I}(\hat{\alpha}_t = 0) - \sum_{t=2}^T \mathbb{I}(\hat{\alpha}_t = \hat{\alpha}_{t-1} \text{ and } \hat{\alpha}_t \hat{\alpha}_{t-1} \neq 0), \\ \text{df}_\lambda(\hat{\beta}) &= \sum_{j=1}^p \left\{ T - \sum_{t=1}^T \mathbb{I}(\hat{\beta}_t^j = 0) - \sum_{t=2}^T \mathbb{I}(\hat{\beta}_t^j = \hat{\beta}_{t-1}^j \text{ and } \hat{\beta}_t^j \hat{\beta}_{t-1}^j \neq 0) \right\}, \\ \text{df}_\lambda(\hat{\gamma}) &= \sum_{j=1}^p \left\{ T - \sum_{t=1}^T \mathbb{I}(\hat{\gamma}_t^j = 0) - \sum_{t=2}^T \mathbb{I}(\hat{\gamma}_t^j = \hat{\gamma}_{t-1}^j \text{ and } \hat{\gamma}_t^j \hat{\gamma}_{t-1}^j \neq 0) \right\}, \end{aligned}$$

are numbers of nonzero and time-varying estimated values in $\hat{\alpha} = (\hat{\alpha}_1, \dots, \hat{\alpha}_T)^T$, $\hat{\beta} = (\hat{\beta}_1^T, \dots, \hat{\beta}_T^T)^T$, and $\hat{\gamma} = (\hat{\gamma}_1^T, \dots, \hat{\gamma}_T^T)^T$, respectively.

Remark 3. There are three regularization parameters $(\lambda_1, \lambda_2, \lambda_3)$ that need to be selected during computation, and three-dimensional searches may be computationally costly. In practical applications, we can adjust for confounders \mathbf{Z}_i in models (1) and (2) during the data preprocessing stage (Schuurmans et al. 2024), and then use the residuals for subsequent analyses. This approach removes the need to include parameters δ and Γ in the model, thereby eliminating the necessity of selecting λ_3 . For the remaining parameters (λ_1, λ_2) , we perform a two-dimensional grid search to identify the combinations which minimize the criterion (16).

6 Simulations

In this section, we conduct comprehensive simulation studies to compare the performance of proposed method with some existing approaches. We consider the following time-varying structural equation models:

$$Y_{it} = \alpha_t X_{it} + \beta_t^T \mathbf{M}_{it} + e_{it}, \quad (17)$$

$$\mathbf{M}_{it} = \gamma_t X_{it} + \boldsymbol{\eta}_t Y_{i,t-1} + \boldsymbol{\varepsilon}_{it}, \quad (18)$$

with $q = 0$ and $T = 4$. We generate the exposure $X_{i1} \sim \text{Uniform}(0, 1)$ and $X_{it} = X_{i,t-1} + \text{Uniform}(0, 1)$ for $t = 2, \dots, T$. The random error vector $(e_{i1}, \dots, e_{iT})^T$ is generated from the multivariate normal distribution $\mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}(\rho_1))$, where $\boldsymbol{\Sigma}(\rho_1)$ is the exchangeable correlation structure with parameter ρ_1 . The random error vector $(\boldsymbol{\varepsilon}_{i1}^T, \dots, \boldsymbol{\varepsilon}_{iT}^T)^T = (\varepsilon_{i1}^1, \dots, \varepsilon_{i1}^p, \dots, \varepsilon_{iT}^1, \dots, \varepsilon_{iT}^p)^T \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}(\rho_2))$, where $\boldsymbol{\Sigma}(\rho_2)$ is the (block-wise) compound symmetry correlation structure with parameter ρ_2 . We set $\rho_1 = 0.1$, $\rho_2 = 0.5$ and $\boldsymbol{\eta}_t = \mathbf{0}$.

We consider the following settings with different sample sizes, different dimensions of mediators, different non-monotone missingness structures, and different missingness mechanisms:

Setting I. We consider the missingness structure (A) illustrated in Figure 4. We set $N = 350$, $p = 50$, and the sample sizes in group 1-4 are (50, 100, 100, 100). The true parameters $\boldsymbol{\alpha} = (0.5, 0.5, 1, 1)^T$, β_t^j and γ_t^j for $j = 1, \dots, 14$ and $t = 1, \dots, 4$ are shown in Figure 5. The β_t^j and γ_t^j are both zero for $j = 15, \dots, 50$ and all t .

Setting II. We consider the missingness structure (B) illustrated in Figure 4, where there is no complete case group. We set $N = 375$, $p = 50$, and the sample sizes in group 1-6 are (75, 75, 75, 50, 50, 50). The true parameters $\boldsymbol{\alpha} = (0.5, 0.5, 1, 1)^T$, β_t^j and γ_t^j for $j = 1, \dots, 14$ and $t = 1, \dots, 4$ are shown in Figure 5. The β_t^j and γ_t^j are both zero for $j = 15, \dots, 50$ and all t .

Setting III. We consider the missingness structure (A) illustrated in Figure 4. We set $N = 550$, $p = 200$, and the sample sizes in group 1-4 are (100, 150, 150, 150). The true parameters $\boldsymbol{\alpha} = (1, 1, 1, 1)^T$, β_t^j and γ_t^j for $j = 1, \dots, 22$ and $t = 1, \dots, 4$ are shown in Figure 6. The β_t^j and γ_t^j are both zero for $j = 23, \dots, 200$ and all t .

Setting IV. We consider the missingness structure (B) illustrated in Figure 4, where there is no complete case group. We set $N = 675$, $p = 200$, and the sample sizes in group 1-6 are (125, 125, 125, 100, 100, 100). The true parameters $\boldsymbol{\alpha} = (1, 1, 1, 1)^T$, β_t^j and γ_t^j for

$j = 1, \dots, 22$ and $t = 1, \dots, 4$ are shown in Figure 6. The β_t^j and γ_t^j are both zero for $j = 23, \dots, 200$ and all t .

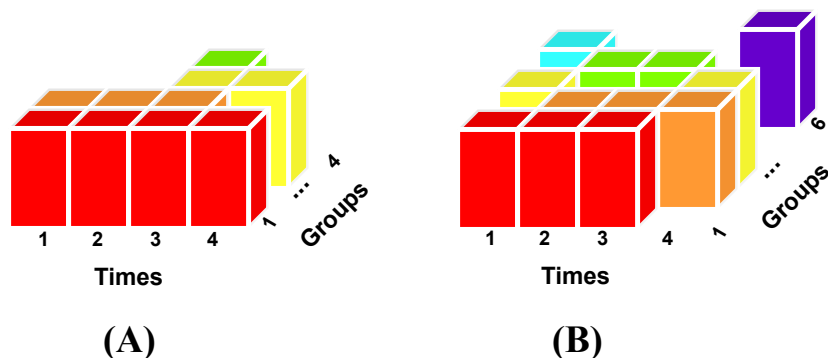


Figure 4: Missingness structures in simulations. Each blank column represents unobserved mediators, while the colored ones represent observed mediators.

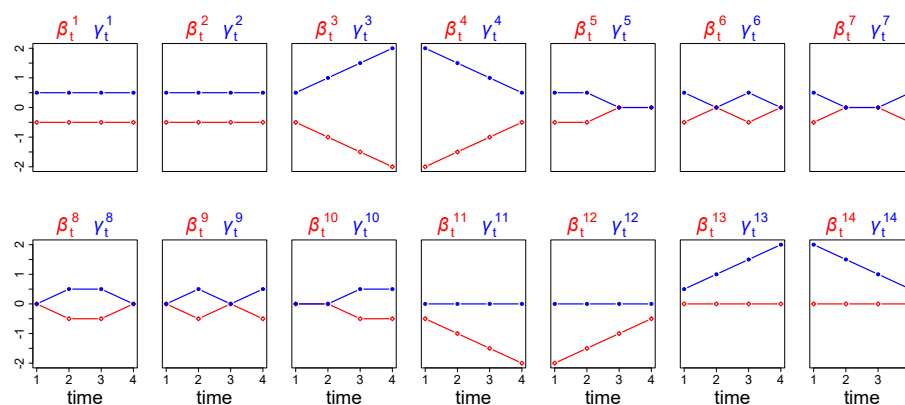


Figure 5: True parameters of β_t^j (red colored) and γ_t^j (blue colored) for $j = 1, \dots, 14$ and $t = 1, \dots, 4$ in simulation settings I and II. The β_t^j and γ_t^j are both zero for $j = 15, \dots, 50$ and all t .

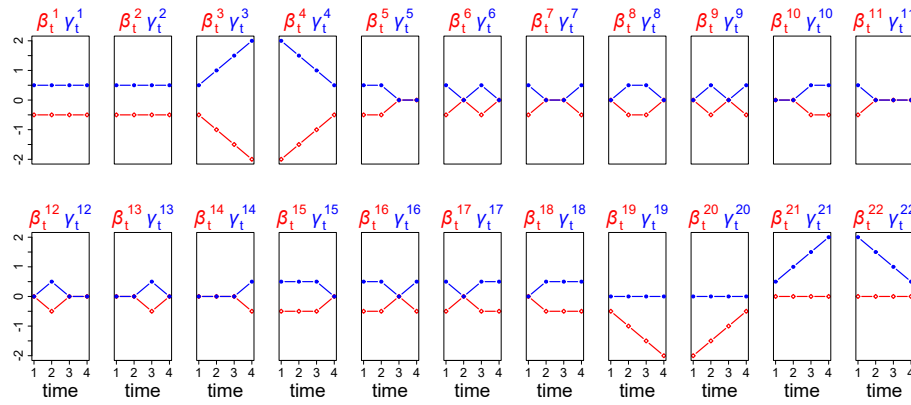


Figure 6: True parameters of β_t^j (red colored) and γ_t^j (blue colored) for $j = 1, \dots, 22$ and $t = 1, \dots, 4$ in simulation settings III and IV. The β_t^j and γ_t^j are both zero for $j = 23, \dots, 200$ and all t .

In each of the settings I-IV, we consider three different missingness mechanisms (Little and Rubin 2019): missing completely at random (MCAR), missing at random (MAR), and missing not at random (MNAR). In the MCAR scenario, all subjects are completely randomly assigned to the groups. In the MAR scenario, subjects $i = 1, \dots, N$ are sequentially randomly assigned to group 1 with probability proportional to $\sum_{t=1}^T |X_{it}|$. The remaining subjects $i \notin \mathcal{H}(1)$ are sequentially randomly assigned to group 2 with probability proportional to $\sum_{t=1}^T |X_{it}|$, and so on until all subjects are grouped. In the MNAR scenario, the grouping process is the same as that in the MAR scenario, except that the probability is proportional to $\sum_{j=1}^p \sum_{t=1}^T |M_{it}^j|$.

Our method, using varying-coefficient (VC) models (17) and (18) with LMI-data, is compared with other mediator selection approaches, including: (i) Mix: mixed-effects models with time interactions and lasso penalty (Rijnhart et al. 2022); (ii) Path: linear structural equation models with pathway lasso penalty (Zhao and Luo 2022); (iii) HIMA: linear structural equation models with de-biased lasso and false discovery rate control (Perera et al. 2022); (iv) Bayes: Bayesian sparse models with continuous shrinkage (Song et al. 2020). Note that for VC and Mix, the parameters are estimated using data from all time points simultaneously, while for Path, HIMA, and Bayes, the parameters are estimated using data from each time point separately. We also consider different approaches for handling missing data, including: (i) CC: complete-case analysis which retains only subjects in group 1 with fully observed repeated measurements; (ii) MC: matrix completion via iterative soft-thresholded singular value decomposition (Mazumder, Hastie and Tibshirani 2010); (iii)

SI: single imputation based on trajectory means (Jahangiri et al. 2023).

In each replication, to evaluate the accuracy of the mediator selection for each method, we calculate the false negative rate (FNR) $\frac{1}{T} \sum_{t=1}^T \frac{\sum_{j=1}^p \mathbb{I}(\beta_t^j \gamma_t^j \neq 0, \hat{\beta}_t^j \hat{\gamma}_t^j = 0)}{\sum_{j=1}^p \mathbb{I}(\beta_t^j \gamma_t^j \neq 0)}$ and the false positive rate (FPR) $\frac{1}{T} \sum_{t=1}^T \frac{\sum_{j=1}^p \mathbb{I}(\beta_t^j \gamma_t^j = 0, \hat{\beta}_t^j \hat{\gamma}_t^j \neq 0)}{\sum_{j=1}^p \mathbb{I}(\beta_t^j \gamma_t^j = 0)}$, across all time points, where $\hat{\beta}_t^j$ and $\hat{\gamma}_t^j$ are estimated values of β_t^j and γ_t^j for $t = 1, \dots, T$ and $j = 1, \dots, p$. We also calculate the mean squared error (MSE) $\frac{1}{T} \sum_{t=1}^T \sum_{j=1}^p (\beta_t^j \gamma_t^j - \hat{\beta}_t^j \hat{\gamma}_t^j)^2$ to evaluate the precision of the mediation effect estimation for each method. We say that a method has better performance if it has a smaller FNR+FPR and a smaller MSE.

Tables 1 and 2 summarize the averaged estimation results across 50 replications of different methods for different missingness mechanisms under settings I and II. Tables 1 shows that the proposed VC-LMI method has the smallest FNR+FPR and MSE among all missingness mechanisms, which suggests that the proposed VC-LMI method has the highest accuracy of mediator selection and precision of mediation effect estimation, and its performance is robust to different missingness mechanisms. The VC-CC has a relatively large FNR+FPR, due to the fact that only 50/350 of the subjects with complete repeated measurements are used for analysis. The VC-MC has large FNR+FPR and MSE, possibly because the MC assumes the low-rank structure and noisy entries of the variable matrix, which may not be suitable for our longitudinal data. The VC-SI also has large FNR+FPR and MSE, possibly because the SI utilizes the trajectory mean for imputation, which erases some of the longitudinal changes in variables. Our VC-LMI method makes full use of the correlation among repeated measurements for imputation, and effectively integrates the multiple imputations, which therefore exhibits powerful performance compared with candidate approaches.

For other mediator selection methods in Tables 1, both FNR+FPR and MSE are large. Specifically, the Mix shows large FPR, while the Path provides large FNR. The reason for their poor estimations may be because the mediator selection methods with some shortcomings (too many interactions in Mix and over-shrinkage in Path) are used on the datasets with poor quality. The HIMA and Bayes provide a smaller FPR but a larger FNR, possibly due to the multiple-testing procedure for HIMA and the posterior inclusion probability for Bayes being overly conservative.

Table 1: False negative rate (FNR), false positive rate (FPR), FNR+FPR, and mean squared error (MSE) of different methods for different missingness mechanisms under setting I.

| | missing completely at random | | | | missing at random | | | | missing not at random | | | |
|---------------|------------------------------|-------|--------------|--------------|-------------------|-------|--------------|--------------|-----------------------|-------|--------------|--------------|
| | FNR | FPR | FNR+FPR | MSE | FNR | FPR | FNR+FPR | MSE | FNR | FPR | FNR+FPR | MSE |
| VC-LMI | 0.059 | 0.086 | 0.145 | 0.719 | 0.049 | 0.090 | 0.139 | 0.774 | 0.064 | 0.088 | 0.152 | 0.715 |
| VC-CC | 0.174 | 0.054 | 0.228 | 1.032 | 0.180 | 0.057 | 0.237 | 1.018 | 0.150 | 0.058 | 0.208 | 0.860 |
| VC-MC | 0.419 | 0.072 | 0.491 | 8.332 | 0.333 | 0.074 | 0.407 | 8.590 | 0.401 | 0.071 | 0.471 | 8.176 |
| VC-SI | 0.203 | 0.036 | 0.239 | 5.065 | 0.180 | 0.034 | 0.214 | 5.228 | 0.196 | 0.036 | 0.232 | 5.208 |
| Mix-CC | 0.111 | 0.105 | 0.215 | 2.638 | 0.109 | 0.114 | 0.222 | 2.525 | 0.116 | 0.099 | 0.215 | 2.319 |
| Mix-MC | 0.086 | 0.327 | 0.414 | 8.815 | 0.077 | 0.320 | 0.397 | 9.214 | 0.092 | 0.314 | 0.407 | 8.772 |
| Mix-SI | 0.057 | 0.167 | 0.224 | 5.839 | 0.043 | 0.163 | 0.205 | 6.016 | 0.049 | 0.157 | 0.206 | 5.965 |
| Path-CC | 0.272 | 0.065 | 0.337 | 7.215 | 0.288 | 0.064 | 0.351 | 7.125 | 0.258 | 0.065 | 0.323 | 7.082 |
| Path-MC | 0.560 | 0.032 | 0.592 | 10.420 | 0.470 | 0.028 | 0.498 | 10.364 | 0.541 | 0.033 | 0.574 | 10.382 |
| Path-SI | 0.248 | 0.020 | 0.268 | 9.198 | 0.227 | 0.022 | 0.250 | 9.319 | 0.249 | 0.021 | 0.270 | 9.285 |
| HIMA-CC | 0.679 | 0.005 | 0.685 | 1.715 | 0.702 | 0.004 | 0.707 | 2.198 | 0.669 | 0.004 | 0.673 | 1.708 |
| HIMA-MC | 0.574 | 0.035 | 0.609 | 8.742 | 0.489 | 0.035 | 0.524 | 9.173 | 0.547 | 0.040 | 0.587 | 8.579 |
| HIMA-SI | 0.376 | 0.014 | 0.390 | 4.895 | 0.361 | 0.015 | 0.377 | 5.066 | 0.370 | 0.015 | 0.385 | 5.158 |
| Bayes-CC | 0.916 | 0.000 | 0.916 | 5.577 | 0.916 | 0.000 | 0.916 | 5.508 | 0.901 | 0.000 | 0.901 | 4.797 |
| Bayes-MC | 0.786 | 0.019 | 0.805 | 10.043 | 0.754 | 0.021 | 0.775 | 10.840 | 0.760 | 0.020 | 0.780 | 9.711 |
| Bayes-SI | 0.689 | 0.001 | 0.690 | 6.129 | 0.675 | 0.001 | 0.676 | 6.129 | 0.659 | 0.001 | 0.660 | 6.115 |

Note: Each estimator is denoted as “mediator selection method - missing data handling method”. Results are averaged across all time points and 50 replications.

Table 2 summarizes the estimation results under setting II corresponding to the missingness structure (B) in Figure 4. Since no subject has complete repeated measurements, the CC is not applicable. Our proposed VC-LMI method has the smallest FNR+FPR and MSE among all candidate approaches. The estimation results under settings III and IV in the high-dimensional context are provided in Tables S1 and S2 in the Supplementary Material, which show that our proposed VC-LMI method also beats other candidate approaches, based on its smallest FNR+FPR and MSE. The Supplementary Material includes additional simulations, where Tables S3 and S4 demonstrate that our method continues to outperform other candidate approaches in higher-dimensional settings ($p = 2000$) and scenarios with more time points ($T = 5$), respectively. Table S5 shows that combining multiple imputations through the GMM outperforms approach that selecting imputations with the smallest variance.

Table 2: False negative rate (FNR), false positive rate (FPR), FNR+FPR, and mean squared error (MSE) of different methods for different missingness mechanisms under setting II.

| | missing completely at random | | | | missing at random | | | | missing not at random | | | |
|---------------|------------------------------|-------|--------------|--------------|-------------------|-------|--------------|--------------|-----------------------|-------|--------------|--------------|
| | FNR | FPR | FNR+FPR | MSE | FNR | FPR | FNR+FPR | MSE | FNR | FPR | FNR+FPR | MSE |
| VC-LMI | 0.046 | 0.105 | 0.151 | 0.790 | 0.048 | 0.113 | 0.160 | 0.855 | 0.056 | 0.107 | 0.163 | 0.755 |
| VC-MC | 0.439 | 0.076 | 0.514 | 9.195 | 0.369 | 0.070 | 0.439 | 9.082 | 0.426 | 0.072 | 0.498 | 8.962 |
| VC-SI | 0.254 | 0.043 | 0.297 | 5.426 | 0.276 | 0.043 | 0.319 | 5.477 | 0.246 | 0.043 | 0.290 | 5.374 |
| Mix-MC | 0.105 | 0.330 | 0.435 | 10.002 | 0.091 | 0.291 | 0.382 | 10.131 | 0.114 | 0.297 | 0.411 | 9.880 |
| Mix-SI | 0.057 | 0.140 | 0.198 | 6.707 | 0.060 | 0.143 | 0.203 | 6.606 | 0.056 | 0.151 | 0.207 | 6.581 |
| Path-MC | 0.646 | 0.022 | 0.668 | 10.908 | 0.559 | 0.017 | 0.577 | 10.770 | 0.611 | 0.018 | 0.630 | 10.851 |
| Path-SI | 0.286 | 0.019 | 0.305 | 9.609 | 0.276 | 0.016 | 0.292 | 9.534 | 0.264 | 0.019 | 0.283 | 9.544 |
| HIMA-MC | 0.654 | 0.045 | 0.698 | 9.667 | 0.571 | 0.046 | 0.618 | 9.837 | 0.631 | 0.048 | 0.680 | 9.564 |
| HIMA-SI | 0.457 | 0.021 | 0.478 | 5.321 | 0.461 | 0.018 | 0.479 | 5.287 | 0.429 | 0.021 | 0.450 | 5.360 |
| Bayes-MC | 0.915 | 0.028 | 0.943 | 10.791 | 0.876 | 0.031 | 0.907 | 11.405 | 0.902 | 0.030 | 0.933 | 10.800 |
| Bayes-SI | 0.756 | 0.008 | 0.764 | 5.960 | 0.774 | 0.005 | 0.780 | 6.075 | 0.764 | 0.006 | 0.770 | 6.010 |

Note: Each estimator is denoted as “mediator selection method - missing data handling method”. Results are averaged across all time points and 50 replications.

7 DNHS data analysis

In this section, we investigate how DNAm dynamically mediates the effect of traumatic experiences on development of PTSD, utilizing the DNHS data. Specifically, we treat the number of traumatic events as exposure X_{it} , the logarithm transformed PTS symptom severity score as outcome Y_{it} , and the M-value transformed (Kruppa et al. 2021) DNAm CpG sites as potential mediators \mathbf{M}_{it} . The study includes 526 subjects in waves 1, 2, 4, and 5, which are referred to as time points 1-4. In our data, exposure X_{it} and outcome Y_{it} are always observed, but the mediators \mathbf{M}_{it} suffer from non-monotone missingness as shown in Figure 1: \mathbf{M}_i of some subjects may be available at certain visits, but missing at the next time point, and measured again at later visits.

There are two steps in the data preprocessing phase. First, we adjust some demographic information and blood work measured at baseline. For the outcome corresponding to model (1), we adjust age, gender, race, education, income, cigarette use, depression, and social cohesion score (Johns et al. 2012, Wani et al. 2021). For the mediators corresponding to model (2), we adjust age, gender, cigarette use, CD4 T cells, CD8 T cells, natural killer cells, B cells, and monocytes (Occean et al. 2022). The residuals after adjustment and the

models (17) and (18) are used for subsequent analyses. Second, given the limited number of subjects and the ultra-high dimensional potential mediators (1879 DNAm CpG sites), we carry out a screening process to reduce the dimension to a moderate scale below the sample size (Fan and Lv 2008). Specifically, for $t = 1, \dots, T$ and $j = 1, \dots, p$, we consider a series of marginal models $Y_{it} = \alpha_t X_{it} + \beta_t^j M_{it}^j + e_{it}$ and $M_{it}^j = \gamma_t^j X_{it} + \eta_t^j Y_{i,t-1} + \varepsilon_{it}^j$. Along the lines of the sure independence mediator screening (Perera et al. 2022), for each time point $t = 1, \dots, T$, we use the subjects with observed mediators to fit the marginal models, then obtain a subset $\mathcal{D}_t = \{j: M^j \text{ is among the top } [n_t/\log(n_t)] \text{ largest } \beta_t^j \gamma_t^j \text{ effect}\}$, where n_t is the sample size at time point t and $[\cdot]$ denotes the integer part. Then we take $\bigcup_{t=1, \dots, T} \mathcal{D}_t$ as the index set of potential mediators after screening. There are 153 DNAm CpG sites remaining for the next step of analysis.

We use the proposed VC-LMI method to fit the preprocessed data. The estimated direct effect is $(\hat{\alpha}_1, \hat{\alpha}_2, \hat{\alpha}_3, \hat{\alpha}_4) = (0.50, 0.54, 0.54, 0.54)$, which means that trauma exposure has a sustained direct effect on PTSD. Figures 7 and 8 show the 13 selected CpG sites, each of which has nonzero $\hat{\beta}_t^j \hat{\gamma}_t^j$ at least at one time point. We can find that certain CpG sites (e.g., cg22564046) initially have nonzero mediation effects, but the effects disappear over time. Certain CpG sites (e.g., cg17318247) do not show mediation effects in the early stages, but the effects emerge over time. This indicates that the effects at the CpG level may vary in different stages of the disease process.

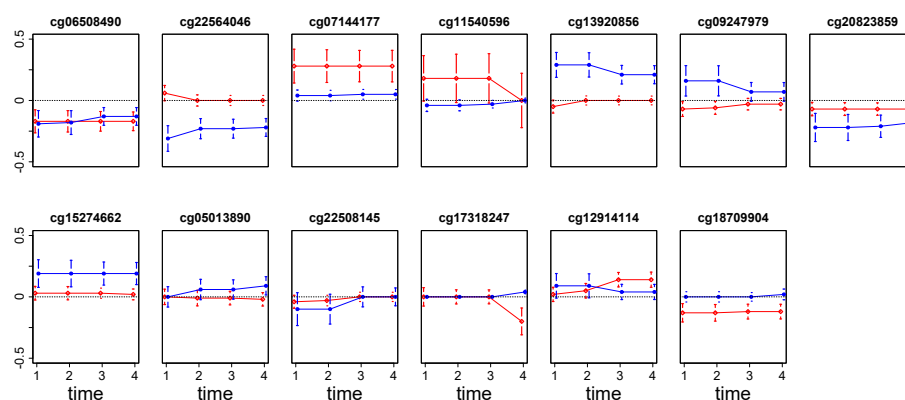


Figure 7: Estimated values $\hat{\beta}_t^j$ (red colored) and $\hat{\gamma}_t^j$ (blue colored) with Bootstrap 95% confidence intervals for the selected CpG sites and $t = 1, \dots, 4$, based on the proposed VC-LMI method.

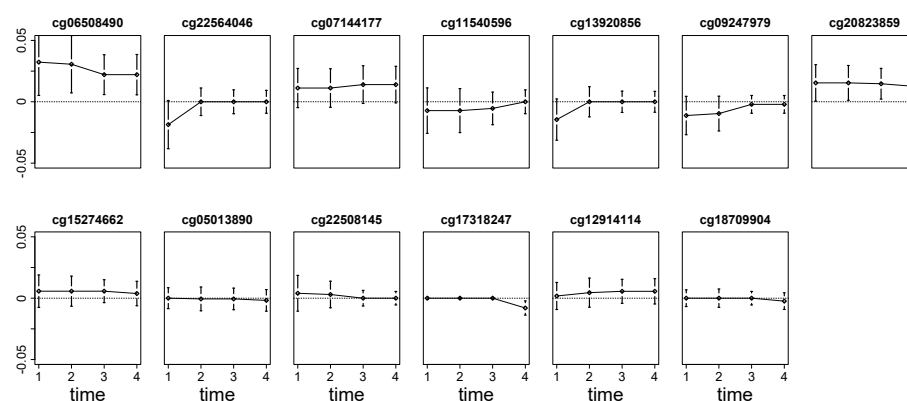


Figure 8: Estimated mediation effects $\hat{\beta}_t^j \hat{\gamma}_t^j$ with Bootstrap 95% confidence intervals for the selected CpG sites and $t = 1, \dots, 4$, based on the proposed VC-LMI method.

As shown in Figure 8, the estimated mediation effects for certain CpG sites at specific time points, such as cg22508145 at $t = 1, 2$ and cg09247979 at $t = 3, 4$, appear small or close to zero. These seemingly negligible effects should be interpreted within the broader context of longitudinal mediation analysis and the cumulative role of DNAm across multiple time points. First, mediation effects can vary over time due to the dynamic nature of DNAm and its interactions with environmental exposures and disease progression. Small mediation effects at individual time points do not necessarily indicate that the CpG site lacks relevance as a mediator. Instead, they may reflect temporal fluctuations or periods during which the indirect effect is less pronounced. Second, while mediation effects may appear negligible at specific time points, the cumulative contribution of a CpG site across all time points could still be meaningful. Longitudinal analyses are designed to capture these aggregate patterns, which may not be evident when focusing solely on isolated single-time-point effects. Finally, even small mediation effects can hold biological significance, especially in high-dimensional settings where individual mediators contribute subtly but collectively to disease mechanisms (Xue et al. 2022).

Table 3 lists the gene information corresponding to the selected CpG sites. Eleven of these sites are annotated to human genes according to the Illumina annotation files (*CAT*, *PFKP*, *EEF1E1*, *TANC1*, *PTPRK*, *OVGP1*, *RPS6KA2*, *CPAMD8*, *CDK16*, *FAM120B*, and *C14orf182*). Of these eleven genes, *PTPRK* has been shown to be associated with mental health disorders such as PTSD (Chitralla, Nagarkatti and Nagarkatti 2016), *RPS6KA2* has been reported as significant KEGG pathways for PTSD in epigenome-wide DNA methylation study (Kuan et al. 2017). Other CpG sites while not located near any protein-coding

genes, fall within ENCODE candidate Cis-Regulatory Elements that are adjacent to highly conserved non-coding genomic regions (Moore et al. 2020).

Table 3: Gene information corresponding to the selected CpG sites.

| CpG | chromosome | Gene |
|------------|------------|-----------|
| cg06508490 | 11 | CAT |
| cg22564046 | 10 | PFKP |
| cg07144177 | 6 | EEF1E1 |
| cg11540596 | 6 | NA/Desert |
| cg13920856 | 2 | TANC1 |
| cg09247979 | 6 | PTPRK |
| cg20823859 | 1 | OVGP1 |
| cg15274662 | 7 | NA/Desert |
| cg05013890 | 6 | RPS6KA2 |
| cg22508145 | 19 | CPAMD8 |
| cg17318247 | X | CDK16 |
| cg12914114 | 6 | FAM120B |
| cg18709904 | 14 | C14orf182 |

Beyond the time-invariant associations, our results further show that the effects of genes change over time, which may reveal how the biological embedding of trauma exposure through DNAm contributes to PTSD risk over development. Our findings could have important clinical implications for time-sensitive risk prediction in trauma-exposed populations.

To compare the performance of the proposed VC-LMI method with other approaches, we randomly split the preprocessed data into a training set (80% of subjects) and a testing set (20% of subjects; the index set is denoted as \mathfrak{T}) for 50 replications. In each replication, for each method, we train the model on the training set and obtain the estimated values of parameters. Then we calculate the mean number of selected mediators (No.-mediator) as $\frac{1}{T} \sum_{t=1}^T \sum_{j=1}^p \mathbb{I}(\hat{\beta}_t^j \hat{\gamma}_t^j \neq 0)$, the prediction mean squared error (PMSE- β) for the outcome model (17) as $\frac{1}{|\mathfrak{T}|} \sum_{i \in \mathfrak{T}} \sum_{r=1}^R \mathbb{I}(i \in \mathcal{H}(r)) \frac{1}{|o(r)|} \sum_{t \in o(r)} (Y_{it} - \hat{Y}_{it})^2$, and the prediction mean squared error (PMSE- γ) for the mediator model (18) as $\frac{1}{|\mathfrak{T}|p} \sum_{i \in \mathfrak{T}} \sum_{r=1}^R \mathbb{I}(i \in \mathcal{H}(r)) \frac{1}{|o(r)|} \sum_{t \in o(r)} (\mathbf{M}_{it} - \hat{\mathbf{M}}_{it})^T (\mathbf{M}_{it} - \hat{\mathbf{M}}_{it})$, on the testing set across

all time points, where \hat{Y}_{it} and $\hat{\mathbf{M}}_{it}$ are the fitted values of Y_{it} and \mathbf{M}_{it} for $i \in \mathcal{I}$ and $t \in o(r)$.

Table 4 summarizes the averaged estimation results across 50 replications. We observe that our VC-LMI method has the smallest PMSE both in the outcome model and the mediator model, indicating that the proposed method is more accurate in terms of prediction. For the Mix and the Path, despite more mediators being selected based on the imputed dataset, the PMSE is larger, possibly due to overfitting. For the HIMA and the Bayesian approaches, no mediator passes their overly conservative significance tests, which is consistent with the simulation results in Section 6.

Table 4: Mean number of selected mediators (No.-mediator), prediction mean squared error (PMSE- β) for the outcome model (17), and prediction mean squared error (PMSE- γ) for the mediator model (18) in DNHS data analysis.

| | No.-mediator | PMSE- β | PMSE- γ |
|---------------|--------------|---------------|----------------|
| VC-LMI | 10.350 | 0.087 | 1.240 |
| VC-CC | 8.900 | 0.286 | 1.754 |
| VC-MC | 14.500 | 0.097 | 1.249 |
| VC-SI | 22.250 | 0.104 | 1.254 |
| Mix-CC | 2.400 | 0.133 | 1.396 |
| Mix-MC | 6.350 | 0.101 | 1.249 |
| Mix-SI | 32.200 | 0.097 | 1.251 |
| Path-CC | 12.300 | 0.125 | 1.385 |
| Path-MC | 43.900 | 0.077 | 1.248 |
| Path-SI | 45.700 | 0.092 | 1.251 |
| HIMA-CC | 0.000 | 0.129 | 1.250 |
| HIMA-MC | 0.150 | 0.111 | 1.250 |
| HIMA-SI | 0.000 | 0.113 | 1.250 |
| Bayes-CC | 0.000 | 0.110 | 1.250 |
| Bayes-MC | 0.000 | 0.113 | 1.250 |
| Bayes-SI | 0.000 | 0.113 | 1.251 |

Note: Each estimator is denoted as “mediator selection method - missing data handling method”. Results are averaged across all time points and 50 replications.

8 Discussion

In this paper, to investigate the time-varying mediation effects of DNAm on the relationship between trauma and PTSD, we propose time-varying structural equation models. The LMI approach is proposed to handle non-monotone missing DNAm and the regularization approach is incorporated to select relevant mediators and capture time-varying effects. In simulations, the proposed method compares favorably against existing competitors in various scenarios. In DNHS data analysis, we successfully identify potential DNAm CpG sites which show dynamic mediation effects.

We assume the mediators ($M_{it}^1, \dots, M_{it}^p$) are correlated but we do not specify causal ordering. The main reason for this assumption is that, based on the existing literature, external adverse environmental exposures or disease conditions are likely to influence DNAm across multiple CpG sites (Turecki and Meaney 2016), whereas the relationships between DNAm levels across CpG sites are primarily correlational (Mou et al. 2022). Such assumptions of parallel relationships among mediators are commonly made in high-dimensional DNAm-related mediation analysis (e.g., Song et al. 2020, Perera et al. 2022, Xue et al. 2022). Furthermore, even if causal ordering existed among high-dimensional DNAm CpG sites, establishing the correct causal order is infeasible in complex biological systems, instead they are characterized by intricate interactions and feedback loops (Tai et al. 2022). In our study, we also lack sufficient genetic knowledge to reliably support causal ordering. A more refined understanding of the causal relationships among high-dimensional mediators may require the application of directed acyclic graph learning techniques (Shi and Li 2022), which might be able to identify the complex causal structures among variables. However, such analyses extend beyond the scope of this study and could be explored in future research.

Unmeasured confounding is an inherent challenge in population-based studies. In our DNHS data analysis, we incorporate a comprehensive set of demographic information and blood work measured at baseline to control for confounding bias as much as possible. Nonetheless, some unmeasured confounders, such as participants' other medical treatments or psychological interventions for PTSD not captured in the DNHS investigation, may still exist. To address these remaining unmeasured confounders, methods such as instrumental variable approaches (Chen et al. 2023), proxy causal learning frameworks (Dukes et al. 2023), or de-confounding techniques (Yuan and Qu 2024) could be employed. Exploring

these methodologies could be a promising direction for future research.

In our imputation step, we estimate $\mathbb{E}(M_t^j | \mathbf{M}_t, \mathbf{D})$ using specific groups that include observed M_t^j and \mathbf{M}_t , then impute the unobserved M_t^j in other groups that contain observed \mathbf{M}_t but missing M_t^j . The rationale for our imputation method relies on the premise that the conditional distributions of mediators in the observed data are the same as those in the missing data, which holds only under the assumptions of MCAR or MAR. For mediators that are MNAR, the direct and indirect effects corresponding to the coefficients in models (1) and (2) are not identifiable without additional assumptions. Potential strategies include modeling the missing indicator $f_{\mathbf{R}_i}(\mathbf{r}_i | \mathbf{X}_i, \mathbf{Y}_i, \mathbf{Z}_i, \mathbf{M}_i)$ based on the completeness assumption (Zuo et al. 2024), or introducing instrumental or shadow variables to impose specific structures on the missingness mechanism (Shan et al. 2024). Extending our method to accommodate MNAR mediators using these strategies, as well as further addressing missing outcomes (Qin et al. 2019), are important directions for future research.

For the statistical inference of the mediation effect, there are two main challenges including the limiting distributions of the estimators for parameters in models (1) and (2), and the composite p -value of the mediation effect. Specifically, the penalty functions in (9) and (10) may present challenges in deriving the asymptotic distribution. Some existing studies used the bi-level penalty similar to that of $P_1(\boldsymbol{\Theta})$ in (9) and the corresponding theoretical results show that the estimator of nonzero coefficients in general converges to nonnormal distribution (Huang et al. 2009). Our estimator includes an additional fusion penalty $P_2(\boldsymbol{\Theta})$ in (10) to serve the purpose of shrinking similar effects at adjacent time points, which may lead to more complex limiting distributions. On the other hand, to test whether M^j mediates the effect of X on Y , the null and alternative hypotheses can be formulated as $H_0^j : \beta^j \gamma^j = 0$ versus $H_1^j : \beta^j \gamma^j \neq 0$. The null hypothesis encompasses a product of parameters, which is composite and consists of three cases as case 1: $\beta^j = 0$ and $\gamma^j \neq 0$, case 2: $\beta^j \neq 0$ and $\gamma^j = 0$, and case 3: $\beta^j = 0$ and $\gamma^j = 0$. Constructing the test statistic and deriving the reference distribution and composite p -value under the composite null hypothesis are still open problems in the field of high-dimensional mediation analysis, and relevant studies can be found in the review by Du et al. (2023). Considering the complex limiting distributions of $\hat{\beta}^j$ and $\hat{\gamma}^j$, it is more challenge to conduct statistical inference for the mediation effect $\hat{\beta}^j \hat{\gamma}^j$. We will develop theories and inference tools of the proposed estimator in our future research.

Funding

The research is supported by NSF Grant DMS-2210640 and 1952406, and in part by NIH Grant R01MD011728.

Supplement

Supplement to “Time-varying mediation analysis for incomplete data with application to DNA methylation study for PTSD”. The Supplementary Material includes three sections. Section 1 contains the GMM estimator corresponding to each mediator M^j ($j = 1, \dots, p$) in model (2). Section 2 contains the expanded version of Algorithm 1. Section 3 contains the additional simulation results.

Supplemental code. We have implemented the proposed longitudinal multiple imputation method as an R function named “LMI” and the generalized method of moments with principal component analysis as another R function named “GMM_PCA”. Simulation code, data, and application code, along with detailed instructions, are provided.

References

- BIND, M. A. C., VANDERWEELE, T. J., COULL, B. A. and SCHWARTZ, J. D. (2016). Causal mediation analysis for longitudinal data with exogenous exposure. *Biostatistics* **17** 122–134.
- BRAUN, P. R., HAN, S., HING, B., NAGAHAMA, Y., GAUL, L. N., HEINZMAN, J. T., GROSSBACH, A. J., CLOSE, L., DLOUHY, B. J., HOWARD III, M. A., KAWASAKI, H., POTASH, J. B. and SHINOZAKI, G. (2019). Genome-wide DNA methylation comparison between live human brain and peripheral tissues within individuals. *Transl. Psychiatry* **9** 47.
- BRESLAU, N., KESSLER, R. C., CHILCOAT, H. D., SCHULTZ, L. R., DAVIS, G. C. and ANDRESKI, P. (1998). Trauma and posttraumatic stress disorder in the community: The 1996 Detroit Area Survey of Trauma. *Arch. gen. psychiatry* **55** 626–632.
- CAI, X., COFFMAN, D. L., PIPER, M. E. and LI, R. (2022). Estimation and inference

- for the mediation effect in a time-varying mediation model. *BMC Medical Res. Methodol.* **22** 113.
- CHEN, C., SHEN, B., LIU, A., WU, R. and WANG, M. (2021). A multiple robust propensity score method for longitudinal analysis with intermittent missing data. *Biometrics* **77** 519–532.
- CHEN, F., HU, W., CAI, J., CHEN, S., SI, A., ZHANG, Y. and LIU, W. (2023). Instrumental variable-based high-dimensional mediation analysis with unmeasured confounders for survival data in the observational epigenetic study. *Front. Genet.* **14** 1092489.
- CHEN, X., LIN, Q., KIM, S., CARBONELL, J. G. and XING, E. P. (2012). Smoothing proximal gradient method for general structured sparse regression. *Ann. Appl. Stat.* **6** 719–752.
- CHITRALA, K. N., NAGARKATTI, P. and NAGARKATTI, M. (2016). Prediction of possible biomarkers and novel pathways conferring risk to post-traumatic stress disorder. *PLOS ONE* **11** e0168404.
- CHO, H. and QU, A. (2015). Efficient estimation for longitudinal data by combining large-dimensional moment conditions. *Electron. J. Stat.* **9** 1315–1334.
- DORAN, H. E. and SCHMIDT, P. (2006). GMM estimators with improved finite sample properties using principal components of the weighting matrix, with an application to the dynamic panel data model. *J. Econom.* **133** 387–409.
- DU, J., ZHOU, X., CLARK-BOUCHER, D., HAO, W., LIU, J., SMITH, J. A. and MUKHERJEE, B. (2023). Methods for large-scale single mediator hypothesis testing: Possible choices and comparisons. *Genet. Epidemiol.* **47** 167–184.
- DUKES, O., SHPITSER, I. and TCHETGEN TCHETGEN, E. J. (2023). Proximal mediation analysis. *Biometrika* **110** 973–987.
- FAN, J. and LI, R. (2001). Variable selection via nonconcave penalized likelihood and its oracle properties. *J. Am. Stat. Assoc.* **96** 1348–1360.
- FAN, J. and LV, J. (2008). Sure independence screening for ultrahigh dimensional feature space. *J. R. Stat. Soc. Ser. B. Stat. Methodol.* **70** 849–911.

- FRIEDMAN, J., HASTIE, T., and TIBSHIRANI, R. (2010). Regularization paths for generalized linear models via coordinate descent. *J. Stat. Softw.* **33** 1–22.
- GE, L., WANG, J., SHI, C., WU, Z. and SONG, R. (2023). A reinforcement learning framework for dynamic mediation analysis. *PMLR* **202** 11050–11097.
- GOLDMANN, E., AIELLO, A. E., UDDIN, M., DELVA, J., KOENEN, K., GANT, L. M. and GALEA, S. (2011). Pervasive exposure to violence and posttraumatic stress disorder in a predominantly African American Urban Community: The Detroit Neighborhood Health Study. *J. Trauma. Stress* **24** 747–751.
- HANSEN, L. P. (1982). Large sample properties of generalized method of moments estimators. *Econometrica* **50** 1029–1054.
- HUANG, J., MA, S., XIE, H. and ZHANG, C. H. (2009). A group bridge approach for variable selection. *Biometrika* **96** 339–355.
- JAHANGIRI, M., KAZEMNEJAD, A., GOLDFELD, K. S., DANESHPOUR, M. S., MOSTAFAEI, S., KHALILI, D., MOGHADAS, M. R. and AKBARZADEH, M. (2023). A wide range of missing imputation approaches in longitudinal data: A simulation study and real data analysis. *BMC Medical Res. Methodol.* **23** 161.
- JOHNS, L. E., AIELLO, A. E., CHENG, C., GALEA, S., KOENEN, K. C. and UDDIN, M. (2012). Neighborhood social cohesion and posttraumatic stress disorder in a community-based sample: Findings from the Detroit Neighborhood Health Study. *Soc. Psychiatry Psychiatr. Epidemiol.* **47** 1899–1906.
- KOENEN, K. C., RATANATHARATHORN, A., NG, L., MCCLAUGHLIN, K. A., BROMET, E. J., STEIN, D. J., KARAM, E. G., RUSCIO, A. M., BENJET, C., SCOTT, K., ATWOLI, L., PETUKHOVA, M., LIM, C. C. W., AGUILAR-GAXIOLA, S., AL-HAMZAWI, A., ALONSO, J., BUNTING, B., CIUTAN, M., DE GIROLAMO, G., DEGENHARDT, L., GUREJE, O., HARO, J. M., HUANG, Y., KAWAKAMI, N., LEE, S., NAVARRO-MATEU, F., PENNELL, B. E., PIAZZA, M., SAMPSON, N., TEN HAVE, M., TORRES, Y., VIANA, M. C., XAVIER, M. and KESSLER, R. C. (2017). Posttraumatic stress disorder in the World Mental Health Surveys. *Psychol. Med.* **47** 2260–2274.

- KRUPPA, J., SIEG, M., RICHTER, G. and POHRT, A. (2021). Estimands in epigenome-wide association studies. *Clin. Epigenet.* **13** 98.
- KUAN, P., WASZCZUK, M. A., KOTOV, R., MARSIT, C. J., GUFFANTI, G., GONZALEZ, A., YANG, X., KOENEN, K., BROMET, E. and LUFT, B. J. (2017). An epigenome-wide DNA methylation study of PTSD and depression in World Trade Center responders. *Transl. Psychiatry* **7** e1158.
- LITTLE, R. J. A. and RUBIN, D. B. (2019). *Statistical Analysis with Missing Data*, 3rd ed. Wiley, New York.
- LUO, L., SHI, C., WANG, J., WU, Z. and LI, L. (2023). Multivariate dynamic mediation analysis under a reinforcement learning framework. *arXiv preprint arXiv:2310.16203*.
- LUSSIER, A. A., ZHU, Y., SMITH, B. J., CERUTTI, J., FISHER, J., MELTON, P. E., WOOD, N. M., COHEN-WOODS, S., HUANG, R. C., MITCHELL, C., SCHNEPER, L., NOTTERMAN, D. A., SIMPKIN, A. J., SMITH, A. D. A. C., SUDERMAN, M. J., WALTON, E., RELTON, C. L., RESSLER, K. and DUNN, E. C. (2023). Association between the timing of childhood adversity and epigenetic patterns across childhood and adolescence: Findings from the Avon Longitudinal Study of Parents and Children (ALSPAC) prospective cohort. *Lancet Child Adolesc. Health* **7** 532–543.
- MAZUMDER, R., HASTIE, T. and TIBSHIRANI, R. (2010). Spectral regularization algorithms for learning large incomplete matrices. *J. Mach. Learn. Res.* **11** 2287–2322.
- MOORE, J. E., PURCARO, M. J., PRATT, H. E., EPSTEIN C. B., SHORESH, N., ADRIAN, J., KAWLI, T., DAVIS, C. A., DOBIN, A., KAUL, R., HALOW, J., VAN NOSTRAND, E. L., FREESE, P., GORKIN DU, SHEN Y, HE Y, MACKIEWICZ M, PAULI-BEHN F, WILLIAMS BA, MORTAZAVI A, KELLER CA, ZHANG XO, ELHAJ-JAJY SI, HUEY J, DICKEL DE, SNETKOVA V, WEI X, WANG X, RIVERA-MULIA JC, ROZOWSKY J, ZHANG J, CHHETRI SB, ZHANG J, VICTORSEN A, WHITE KP, VISEL A, YEO GW, BURGE CB, LÉCUYER E, GILBERT DM, DEKKER J, RINN J, MENDENHALL EM, ECKER JR, KELLIS M, KLEIN RJ, NOBLE WS, KUNDAJE A, GUIGÓ R, FARNHAM PJ, CHERRY JM, MYERS RM, REN B, GRAVELEY BR, GERSTEIN MB, PENNACCHIO LA, SNYDER MP, BERNSTEIN BE, WOLD B, HARDISON

- RC, GINGERAS TR, STAMATOYANNOPOULOS JA and WENG Z (2020). Expanded encyclopaedias of DNA elements in the human and mouse genomes. *Nature* **583** 699–710.
- MOU, X., ZHANG, H. and ARSHAD, S. H. (2022). Identifying intergenerational patterns of correlated methylation sites. *Ann. Appl. Stat.* **16** 521-536.
- OCCEAN, J. R., WANI, A. H., DONGLASAN, J., AIELLO, A. E., GALEA, S., KOENEN, K. C., QU, A., WILDMAN, D. and UDDIN, M. (2022). DNA methylation of Nuclear Factor of Activated T Cells 1 mediates the prospective relation between exposure to different traumatic event types and post-traumatic stress disorder. *Psychiatry Res.* **311** 114510.
- PERERA, C., ZHANG, H., ZHENG, Y., HOU, L., QU, A., ZHENG, C., XIE, K. and LIU, L. (2022). HIMA2: High-dimensional mediation analysis and its application in epigenome-wide DNA methylation data. *BMC Bioinform.* **23** 296.
- QIN, X., HONG, G., DEUTSCH, J. and BEIN, E. (20219). Multisite causal mediation analysis in the presence of complex sample and survey designs and non-random non-response. *J. R. Stat. Soc., A: Stat. Soc.* **182** 1343–1370.
- RIJNHART, J. J. M., TWISK, J. W. R., VALENTE, M. J. and HEYMANS, M. W. (2022). Time lags and time interactions in mixed effects models impacted longitudinal mediation effect estimates. *J. Clin. Epidemiol.* **151** 143–150.
- RUGGIERO, K. J., BEN, K. D., SCOTTI, J. R. and RABALAIS, A. E. (2003). Psychometric properties of the PTSD Checklist - Civilian version. *J. Trauma. Stress* **16** 495–502.
- RUTTEN, B. P. F., VERMETTEN, E., VINKERS, C. H., URSINI, G., DASKALAKIS, N. P., PISHVA, E., DE NIJS, L., HOUTEPEN, L. C., EIJSSEN, L., JAFFE, A. E., KENIS, G., VIECHTBAUER, W., VAN DEN HOVE, D., SCHRAUT, K. G., LESCH, K. P., KLEINMAN, J. E., HYDE, T. M., WEINBERGER, D. R., SCHALKWYK, L., LUNNON, K., MILL, J., COHEN, H., YEHUDA, R., BAKER, D. G., MAIHOFFER, A. X., NIEVERGELT, C. M., GEUZE, E. and BOKS, M. P. M. (2018). Longitudinal analyses of the DNA methylome in deployed military servicemen identify susceptibility loci for post-traumatic stress disorder. *Mol. Psychiatry* **23** 1145–1156.

- SCHUURMANS, I. K., DUNN, E. C. and LUSSIER, A. A. (2024). DNA methylation as a possible causal mechanism linking childhood adversity and health: Results from two-sample mendelian randomization study. *Am. J. Epidemiol.* **193** 1541–1552.
- SHAN, J., LI, W. and AI, C. (2024). Efficient nonparametric inference of causal mediation effects with nonignorable missing confounders. *arXiv preprint arXiv:2402.05384*.
- SHI, C. and LI, L. (2022). Testing mediation effects using logic of boolean matrices. *J. Am. Stat. Assoc.* **117** 2014–2027.
- SONG, Y., ZHOU, X., ZHANG, M., ZHAO, W., LIU, Y., KARDIA, S. L. R., DIEZ ROUX, A. V., NEEDHAM, B. L., SMITH, J. A. and MUKHERJEE, B. (2020). Bayesian shrinkage estimation of high dimensional causal mediation effects in omics studies. *Biometrics* **76** 700–710.
- TAI, A. S., LIN, P. H., HUANG, Y. T. and LIN, S. H. (2022). Path-specific effects in the presence of a survival outcome and causally ordered multiple mediators with application to genomic data. *Stat. Methods Med. Res.* **31** 1916–1933.
- TIBSHIRANI, R., SAUNDERS, M., ROSSET, S., ZHU, J. and KNIGHT, K. (2005). Sparsity and smoothness via the fused lasso. *J. R. Stat. Soc. Ser. B. Stat. Methodol.* **67** 91–108.
- TORTELLA-FELIU, M., FULLANA, M. A., PÉREZ-VIGIL, A., TORRES, X., CHAMORRO, J., LITTARELLI, S. A., SOLANES, A., RAMELLA-CRAVARO, V., VILAR, A., GONZÁLEZ-PARRA, J. A., ANDERO, R., REICHENBERG, A., MATAIX-COLS, D., VIETA, E., FUSAR-POLI, P., IOANNIDIS, J. P. A., STEIN, M. B., RADUA, J. and FERNÁNDEZ DE LA CRUZ, L. (2019). Risk factors for posttraumatic stress disorder: An umbrella review of systematic reviews and meta-analyses. *Neurosci. Biobehav. Rev.* **107** 154–165.
- TSENG, C. H., ELASHOFF, R., LI, N. and LI, G. (2016). Longitudinal data analysis with non-ignorable missing data. *Stat. Methods Med. Res.* **25** 205–220.
- TURECKI, G. and MEANEY, M. J. (2016). Effects of the social environment and stress on glucocorticoid receptor gene methylation: A systematic review. *Biol. Psychiatry* **79** 87–96.

- UDDIN, M., AIELLO, A. E., WILDMAN, D. E., KOENEN, K. C., PAWELEC, G., DE LOS SANTOS, R., GOLDMANN, E. and GALEA, S. (2010). Epigenetic and immune function profiles associated with posttraumatic stress disorder. *PNAS* **107** 9470–9475.
- VANDERWEELE, T. J. and TCHETGEN TCHETGEN, E. J. (2017). Mediation analysis with time varying exposures and mediators. *J. R. Stat. Soc. Ser. B. Stat. Methodol.* **79** 917–938.
- WANI, A. H., AIELLO, A. E., KIM, G. S., XUE, F., MARTIN, C. L., RATANATHARATHORN, A., QU, A., KOENEN, K., GALEA, S., WILDMAN, D. and UDDIN, M. (2021). The impact of psychopathology, social adversity and stress-relevant DNA methylation on prospective risk for post-traumatic stress: A machine learning approach. *J. Affect. Disord.* **282** 894–905.
- WEI, K., ZHU, H., QIN, G., ZHU, Z. and TU, D. (2022). Multiply robust subgroup analysis based on a single-index threshold linear marginal model for longitudinal data with dropouts. *Stat. Med.* **41** 2822–2839.
- WILKER, S., VUKOJEVIC, V., SCHNEIDER, A., PFEIFFER, A., INERLE, S., PAULY, M., ELBERT, T., PAPASSOTIROPOULOS, A., DE QUERVAIN, D. and KOLASSA, I. T. (2023). Epigenetics of traumatic stress: The association of NR3C1 methylation and posttraumatic stress disorder symptom changes in response to narrative exposure therapy. *Transl. Psychiatry* **13** 14.
- WU, S., RENZHO, A. M. N., HALL, B. J., SHI, L., LING, L. and CHEN, W. (2021). Time-varying associations of pre-migration and post-migration stressors in refugees’ mental health during resettlement: A longitudinal study in Australia. *Lancet Psychiatry* **8** 36–47.
- XUE, F. and QU, A. (2021). Integrating multisource block-wise missing data in model selection. *J. Am. Stat. Assoc.* **116** 1914–1927.
- XUE, F., TANG, X., KIM, G., KOENEN, K. C., MARTIN, C. L., GALEA, S., WILDMAN, D., UDDIN, M. and QU, A. (2022). Heterogeneous mediation analysis on epigenomic ptsd and traumatic stress in a predominantly African American cohort. *J. Am. Stat. Assoc.* **117** 1669–1683.

- YUAN, Y. and QU, A. (2024). De-confounding causal inference using latent multiple-mediator pathways. *J. Am. Stat. Assoc.* **119** 2051–2065.
- ZENG, S., LANGE, E. C., ARCHIE, E. A., CAMPOS, F. A., ALBERTS, S. C. and LI, F. (2022). A causal mediation model for longitudinal mediators and survival outcomes with an application to animal behavior. *J. Agric. Biol. Environ. Stat.* **28** 197–218.
- ZENG, S., ROSENBAUM, S., ALBERTS, S. C., ARCHIE, E. A. and LI, F. (2021). Causal mediation analysis for sparse and irregular longitudinal data. *Ann. Appl. Stat.* **15** 747–767.
- ZHAO, Y. and LUO, X. (2022). Pathway lasso: Pathway estimation and selection with high-dimensional mediators. *Stat. Its Interface* **15** 39–50.
- ZUO, S., GHOSH, D., DING, P. and YANG, F. (2024). Mediation analysis with the mediator and outcome missing not at random. *J. Am. Stat. Assoc.* 1–12.