

# 多重并行中介分析方法的比较研究

于洋<sup>1</sup> 仇沁晓<sup>1</sup> 尤东方<sup>1</sup> 赵杨<sup>1,2</sup>

<sup>1</sup>南京医科大学公共卫生学院生物统计学系, 南京 211166; <sup>2</sup>南京医科大学生物医学大数据重点实验室/肿瘤个体化医学协同创新中心, 南京 211166

通信作者: 赵杨, Email: yzhao@njmu.edu.cn

**【摘要】目的** 介绍 4 种多重并行中介模型的分析方法, 包括纯回归法、逆概率加权法、扩展的自然效应模型和基于权重的填补法, 并对其进行探讨和比较。**方法** 针对多重并行中介模型, 通过 3 种情境的模拟试验比较不同方法在不同情境下估计直接效应和间接效应的表现, 并应用英国生物样本库的数据集进行实例分析。**结果** 模拟试验和实例分析结果显示纯回归法和逆概率加权法对各效应的估计偏倚较小, 扩展的自然效应模型次之, 基于权重的填补法与另外 3 种的估计结果差异较大。**结论** 不同的多重并行中介分析方法有不同的适用情境以及各自的优缺点, 纯回归法更适用于连续中介的情形, 逆概率加权法更适用于二分类中介的情形, 扩展的自然效应模型在用于两个并行中介的残差呈正相关且相关程度较小时更佳, 而基于权重的填补法可能并不适用于并行中介的情形, 因而实际应用时应根据具体情境选择合适的方法。

**【关键词】** 多重并行中介分析; 并行; 中介效应; 逆概率加权法

## A comparative study of multiple parallel mediation analysis methods

Yu Yang<sup>1</sup>, Qiu Qinxiao<sup>1</sup>, You Dongfang<sup>1</sup>, Zhao Yang<sup>1,2</sup>

<sup>1</sup>Department of Biostatistics, School of Public Health, Nanjing Medical University, Nanjing 211166, China; <sup>2</sup>Key Laboratory of Biomedical Big Data/Cancer Individualized Medicine Collaborative Innovation Center, Nanjing Medical University, Nanjing 211166, China

Corresponding author: Zhao Yang, Email: yzhao@njmu.edu.cn

**【Abstract】 Objective** To introduce and compare four analysis methods of multiple parallel mediation model, including pure regression method, method based on inverse probability weighting, extended natural effect model method and weight-based imputation strategies. **Methods** For the multiple parallel mediation model, the simulation experiments of three scenarios were carried out to compare the performance of different methods in estimating direct and indirect effects in different situations. Dataset from UK Biobank was then analyzed by using the four methods. **Results** The estimation biases of the regression method and the inverse probability weighting method were relatively small, followed by the extended natural effect model method, and the estimation results of the weight-based imputation strategies were quite different from the other three methods. **Conclusions** Different multiple parallel mediation analysis methods have different application situations and their own advantages and disadvantages. The regression method is more suitable for continuous mediator, and the inverse probability weighting method is more suitable for binary mediator. The extended natural effect model method has better performances when the residuals of two parallel mediators are positively correlated and the correlation degree is small. The weight-based imputation strategies might not be appropriate for parallel mediation analysis. Therefore, appropriate methods should be selected according to the specific situation in practice.

**【Key words】** Multiple parallel mediation analysis; Parallel; Mediating effect; Inverse probability weighting

DOI: 10.3760/cma.j.cn112338-20211022-00814

收稿日期 2021-10-22 本文编辑 万玉立

引用格式: 于洋, 仇沁晓, 尤东方, 等. 多重并行中介分析方法的比较研究[J]. 中华流行病学杂志, 2022, 43(5): 739-746.

DOI: 10.3760/cma.j.cn112338-20211022-00814.

Yu Y, Qiu QX, You DF, et al. A comparative study of multiple parallel mediation analysis methods[J]. Chin J Epidemiol, 2022, 43(5):739-746. DOI: 10.3760/cma.j.cn112338-20211022-00814.



在心理学、经济学、医学等众多领域中,研究者常常需要探讨各种变量之间的关系。自变量除了对因变量直接作用外,也有可能通过一个中间变量间接地对因变量产生作用。这个中间变量被称为中介变量,对应的两种作用分别称为直接效应和间接效应。中介分析的目的主要就是探索自变量X与因变量Y之间的因果关系机制,将自变量X与因变量Y之间的因果路径进行分解,判断中介变量在其因果路径中如何起作用。Baron和Kenny<sup>[1]</sup>首先对简单中介分析模型基于线性回归分析的方法提出了直接和间接效应的参数估计和假设检验方法。自1992年Robins和Greenland<sup>[2]</sup>基于反事实框架下提出了因果中介效应的定义后,中介分析研究有了很大的发展<sup>[3-4]</sup>。

传统的中介分析往往只考虑一个中介因素。近年来,研究者们越来越多地关注于探索自变量与因变量之间存在多个中介变量的情形<sup>[5-6]</sup>。这类中介模型被称为多重中介模型(multiple mediation model)。例如,有研究表明患者的文化因素与30 d再入院间的关联是通过出院实践和护理过渡中介介导的<sup>[7]</sup>。在该研究中,研究者对并行中介和顺序中介进行了联合中介分析。又如,在探讨母亲超重和子代超重之间的关联时,对出生方式和肠道菌群进行了顺序中介分析。研究发现母亲超重会通过影响出生方式,继而影响肠道菌群,最终影响子代超重的结局<sup>[8]</sup>。

本文将介绍多重并行中介分析的4种常用方法,即纯回归法、逆概率加权法、扩展的自然效应模型和基于权重的填补法;在并行中介的情形下,通过模拟研究比较4种方法在不同的情境下,估计直接效应和间接效应的表现。再将这4种方法应用于英国生物样本库(UK Biobank)数据库中水果摄入与心血管疾病关联的中介研究,以了解其实际应用中的表现。最后对这4种多重并行中介分析方法的适用情境以及各自的优缺点进行总结与讨论。

### 一、基本原理

**多重中介效应分析方法:**多重中介分析旨在探究自变量通过多个不同中介变量对因变量产生作用的机制。根据多个中介变量在自变量和因变量间起作用的方式不同,多重中介效应又分为并行和顺序多重中介效应。前者是指多个中介变量同时在自变量和因变量间起作用;而后者指多个中介变量间出现顺序性特征,在自变量和因变量间形成中介链<sup>[9-11]</sup>。多重中介模型相对于简单中介模型不仅

可以得到总的中介效应,还可以研究每个中介的特定中介效应或特定路径效应。

在反事实框架下,假设有多个感兴趣的中介, $M = (M^{(1)}, \dots, M^{(K)})$ 。令 $M_a$ 为暴露 $A=a$ 时,中介变量 $M$ 的反事实取值。令 $Y_{am}$ 为 $A=a, M=m$ 时,结局 $Y$ 的反事实取值。估计自然直接效应(natural direct effect, NDE)时,每个个体的中介变量取值为某一暴露水平下,其可能的自然取值。那么比较 $A=a$ 与 $A=a^*$ ,暴露 $A$ 对结局 $Y$ 的NDE时,使中介 $M$ 为 $A=a^*$ 时其可能的自然取值,即NDE定义为 $Y_{aM_a} - Y_{a^*M_a}$ 。相应地,暴露通过中介对结局产生的自然间接效应(natural indirect effect, NIE)定义为 $Y_{aM_a} - Y_{aM_a^*}$ 。NIE假定暴露设置为某一水平 $A=a$ ,然后比较中介 $M$ 在 $A=a$ 时自然取值产生的结局与中介 $M$ 在 $A=a^*$ 时自然取值产生的结局。为了使结果具有因果关系,必须对混杂因素进行控制。NDE和NIE的识别需要满足4个假设:① $A$ - $Y$ 间没有未测量的混杂,即 $Y_{am} \perp\!\!\!\perp A|C$ ;② $M$ - $Y$ 间没有未测量的混杂,即 $Y_{am} \perp\!\!\!\perp M|A, C$ ;③ $A$ - $M$ 间没有未测量的混杂,即 $M_a \perp\!\!\!\perp A|C$ ;④ $M$ - $Y$ 间没有受 $A$ 影响的混杂,即 $Y_{am} \perp\!\!\!\perp M_a|A, C$ 。

1. 纯回归法:在反事实框架下,Vanderweele和Vansteelandt<sup>[12]</sup>将简单中介模型下Baron和Kenny<sup>[1]</sup>提出的乘积法推广到具有多个中介变量的模型中,建立新的多重中介变量模型。涉及模型包括 $Y$ 模型(结局对暴露、所有中介变量和混杂因素的模型)和 $M$ 模型(各中介变量对暴露和混杂因素的模型)。例如,在并行中介情形下,当结局和中介变量是连续型变量时,模型可表示为:

$$E[Y|A, M, C] = \theta_0 + \theta_1 A + \theta_2^{(1)} M^{(1)} + \theta_2^{(2)} M^{(2)} + \dots + \theta_2^{(K)} M^{(K)} + \theta_4' C$$

$$E[M^{(i)}|A, C] = \beta_0^{(i)} + \beta_1^{(i)} A + \beta_2^{(i)} C \text{ for } i = 1, \dots, K$$

其中, $A$ 表示暴露变量, $M$ 表示中介变量, $C$ 表示协变量, $K$ 为中介变量的个数。

$$\text{NDE 可表示为 } E[Y_{aM_a} - Y_{a^*M_a} | c] = \theta_1(a - a^*),$$

$$\text{NIE 可表示为 } E[Y_{aM_a} - Y_{aM_a^*} | c] = [\beta_1^{(1)}\theta_2^{(1)} + \dots + \beta_1^{(K)}\theta_2^{(K)}](a - a^*).$$

纯回归法也可用于顺序中介的情形。

2. 逆概率加权法:该方法对顺序和并行中介类型不做假设。其核心思想是,通过建立 $A$ 模型(暴露对混杂因素的模型)和 $M$ 模型(各中介变量对暴

露和混杂因素的模型),得到每个个体*i*的权重

$$w_i = \frac{I(A_i = a)}{Pr(A_i = a | C_i = c_i)} \times \frac{Pr(M_{1i} = m_{1i} | A_i = a^*, C_i = c_i)}{Pr(M_{1i} = m_{1i} | A_i = a, C_i = c_i)} \times \frac{Pr(M_{2i} = m_{2i} | A_i = a^*, C_i = c_i)}{Pr(M_{2i} = m_{2i} | A_i = a, C_i = c_i)}$$

再通过纯回归法建立结局*Y*关于暴露、所有中介变量和协变量的模型,同时赋予相应的权重。从而估计NDE为 $E[Y_{aM_c} - Y_{a^*M_c} | c] = \theta_1(a - a^*)$ ,NIE为 $E[Y_{aM_c} - Y_{aM_c^*} | c] = [\beta_1^{(1)}\theta_2^{(1)} + \dots + \beta_1^{(K)}\theta_2^{(K)}](a - a^*)^{13-14}$ 。

3. 扩展的自然效应模型:Lange等<sup>[13]</sup>将简单中介模型下的自然效应模型法扩展到了具有多个并行中介变量的情形中<sup>[15]</sup>。该方法的基本原理:首先根据原始数据集建立暴露和各中介的回归模型;再引入辅助变量 $A^1, \dots, A^K$ ,重复原始数据 $2^K$ 次,构建一个扩展数据集,并计算每一行的权重;最后拟合结局的回归模型(只包含 $A, A^1, \dots, A^K$ ),并加权,以估计特定路径的间接效应。具体估计步骤:①首先利用原始数据集,以混杂因素为条件,估计暴露的合适模型;②再利用原始数据集,以暴露和混杂为条件,对每个中介变量估计一个合适的模型;③假设在以暴露和混杂为条件的情况下,各中介间相互独立。引入新变量 $A^1, \dots, A^K$ 为辅助暴露变量,重复原始数据的每一条观测 $2^K$ 次,构建一个扩展数据集。以二元暴露、两个中介为例,则每个观测被重复 $2^2=4$ 次。先让 $A^1$ 取值1,再取值0,将原始数据集重复两次;接下来,该数据集再重复两次,这次让 $A^2$ 先取值1,再取值0,得到的扩展数据集见表1。存在*K*个中介,该过程就重复*K*次,从而得到最终

扩展数据集;④计算权重 $W_i = \frac{1}{P(A = A_i | C = C_i)} \prod_{k=1}^K \frac{P(M^k = M_i^k | A = A_i^k, C = C_i)}{P(M^k = M_i^k | A = A_i, C = C_i)}$ ,其中*i*指代扩展数据集中的第*i*行;⑤最后对结局变量拟合回归模型(如logistic, Cox等),该模型只包含 $A, A^1, \dots, A^K$ ,并赋予权重 $W_i$ 。模型可表示为 $g\left(E\left[Y_{A, M_1^1, \dots, M_k^k}\right]\right) = \alpha + \beta_0 A + \sum_{k=1}^K \beta_k A^k$ ,模型中也可以包括合适的交互作用项。估计NDE为 $E\left[Y_{aM_c} - Y_{a^*M_c}\right] = \beta_0(a - a^*)$ ,NIE为 $E\left[Y_{aM_c} - Y_{aM_c^*}\right] = [\beta_1 + \beta_2 + \dots + \beta_K](a - a^*)$ 。其中, $a, a^*$ 为变量*A*的两个可能取值。

4. 基于权重的填补法:Steen等<sup>[16]</sup>将Vansteelandt等<sup>[17]</sup>提出的针对单个中介变量的基于权重的填补法扩展到存在多个顺序中介变量的情形,为多个中介变量拟合自然效应模型。该方法可以处理多个不同类型的中介变量以及不同类型的结局的情形。其核心思想是,对反事实结局拟合一个只包含辅助变量( $a, a', a''$ )的回归模型,并进行加权。权重的选择依赖于中介模型的选择。现以一个二分类暴露和两个顺序中介为例,具体估计步骤:①为第一个中介变量拟合一个合适的模型,以 $M_1$ 为二分类为例,  $\text{logit}P(M_1 = 1 | A, C) = \beta_0 + \beta_1 A + \beta_2^T C$ ;或为第二个中介变量拟合一个合适的模型,以 $M_2$ 为连续型变量为例,  $f(M_2 | A, M_1, C) = N(\gamma_0 + \gamma_1 A + \gamma_2 M_1 + \gamma_3 A M_1 + \gamma_4^T C, \sigma^2)$ 。②为结局变量拟合一个合适的模型,以二分类结局为例,  $\text{logit}P(Y = 1 | A, M_1, M_2, C) = \delta_0 + \delta_1 A + \delta_2 M_1 + \delta_3 M_2 + \delta_4 A M_1 + \delta_5 A M_2 + \delta_6 M_1 M_2 + \delta_7 A M_1 M_2 + \gamma_8^T C$ 。

表1 原始数据集中前两行扩展的演示说明(扩展的自然效应模型)

原始数据集					扩展数据集						
编号	暴露(A)	中介(M)	结局(Y)	混杂(C)	编号	暴露(A)	新变量(A <sup>1</sup> )	新变量(A <sup>2</sup> )	中介(M)	结局(Y)	混杂(C)
101	1	1	0	1	101	1	0	0	1	0	1
102	0	1	1	0	101	1	1	0	1	0	1
...	...	...	...	...	101	1	0	1	1	0	1
...	...	...	...	...	101	1	1	1	1	0	1
...	...	...	...	...	102	0	0	0	1	1	0
...	...	...	...	...	102	0	1	0	1	1	0
...	...	...	...	...	102	0	0	1	1	1	0
...	...	...	...	...	102	0	1	1	1	1	0
...	...	...	...	...	...	...	...	...	...	...	...

注:...:数据集中后续个体的记录



③引入新变量  $a, a', a''$  构建一个扩展数据集:对每一个个体  $i$ , 第一次重复  $a$  取观测到的暴露值  $A_i$ , 第二次重复取反事实暴露  $1-A_i$ 。这两次重复  $a'$  和  $a''$  都取值为观测到的暴露水平。接着再重复上述步骤一次:若更确信第一个中介模型是正确的, 则选择第一个中介变量的模型, 则令  $a'$  取反事实暴露水平  $1-A_i, a''$  取实际观测到的暴露水平; 若更确信第二个中介模型是正确的, 则选择第二个中介变量的模型, 则令  $a''$  取反事实暴露水平  $1-A_i, a'$  取实际观测到的暴露水平, 得到的扩展数据集见表 2。④根据所选择的中介模型以及相应得到的扩展数据集计算权重:若确信第一个中介模型是正确的, 则选择第一个中介变量的模型,  $W_{1i, a'} = \frac{P(M_1 = M_{1i} | A = a', C_i)}{P(M_1 = M_{1i} | A = a'', C_i)} = \frac{P(M_1 = M_{1i} | A = a', C_i)}{P(M_1 = M_{1i} | A = A_i, C_i)}$ ; 若确信第二个中介模型是正确的, 则选择第二个中介变量的模型,  $W_{2i, a''} = \frac{f(M_2 = M_{2i} | A = a'', M_{1i}, C_i)}{f(M_2 = M_{2i} | A = a', M_{1i}, C_i)} = \frac{f(M_2 = M_{2i} | A = a'', M_{1i}, C_i)}{f(M_2 = M_{2i} | A = A_i, M_{1i}, C_i)}$ ; 如果是多个顺序中介, 可以选用信息量最大的那个中介模型。⑤在原数据集中拟合第 2 步中的结局模型, 得到拟合值  $\hat{E}(Y_i | A = a, M_{1i}, M_{2i}, C_i)$ , 将其作为扩展数据集中反事实结局  $Y_i(a, M_{1i}(a'), M_{2i}(a'', M_{1i}(a')))$  的取值。⑥在扩展数据集中对填补的反事实结局拟合自然效应模型, 该模型只包含  $a, a', a''$ , 同时赋予第 4 步中的权重。NDE 和 NIE 的估计同扩展的自然效应模型, 最后可利用 Bootstrap 法得到各效应估计的标准误和置信区间。

### 二、模拟研究

1. 目的:本研究在不同模拟情形下模拟包含一个暴露变量  $A$ 、两个中介变量  $M_1, M_2$ 、一个二分类结局  $Y$  和两个协变量  $C_1, C_2$  的数据, 计算并比较上述 4 种多重并行中介分析方法在估计中介效应方面的性能, 为研究者选择合适的中介分析方法提供建议。

2. 模拟研究方案:本研究参考已有的模拟研究<sup>[17-19]</sup>, 设置两个二分类协变量  $C = (C_1, C_2)$ , 且  $P(C_1 = 1) = 0.2, P(C_2 = 1) = 0.5$ ; 一个二分类暴露变量  $A$  且  $\text{logit}P(A = 1 | C) = -1 + 0.8C_1 + 0.3C_2$ ; 两个并行中介变量  $M_1, M_2$ , 当  $M_1, M_2$  均为连续型变量

表 2 原始数据集的扩展演示(基于权重的填补法)

选择中介模型		扩展原始数据集			
$P(M_1   A, C)$	$i$	$A_i$	$a$	$a'$	$a''$
	101	1	1	1	1
	101	1	0	1	1
	101	1	1	0	1
	101	1	0	0	1
	102	0	0	0	0
	102	0	1	0	0
	102	0	0	1	0
	102	0	1	1	0
	...	...	...	...	...
$P(M_2   A, M_1, C)$	$i$	$A_i$	$a$	$a'$	$a''$
	101	1	1	1	1
	101	1	0	1	1
	101	1	1	1	0
	101	1	0	1	0
	102	0	0	0	0
	102	0	1	0	0
	102	0	0	0	1
	102	0	1	0	1
	...	...	...	...	...

时,  $f(M_1 | A, C_1, C_2) = N(1 + \log(IE1) \times A + 0.9C_1 + 0.2C_2 + e_1, 1), f(M_2 | A, C_1, C_2) = N(0.5 + \log(IE2) \times A + 0.7C_1 + 0.2C_2 + e_2, 1)$ , 其中  $e_1, e_2 \sim N(0, 1)$ ; 当  $M_1, M_2$  均为二分类变量时,  $\text{logit}P(M_1 = 1 | A, C) = 1 + \log(IE1) \times A + 0.3C_1 + 0.8C_2, \text{logit}P(M_2 = 1 | A, C) = 0.5 + \log(IE2) \times A + 0.1C_1 + 0.6C_2$ ; 一个二分类结局变量  $Y, \text{logit}P(Y = 1 | A, M_1, M_2, C) = -1.5 + \log(0.2) \times A + M_1 + M_2 + 0.7C_1 + 0.4C_2$ 。模拟示意图见图 1。模拟研究将从以下 3 个情境分别评价这 4 种多重并行中介分析方法(表 3)。所有情形的直接效应均设置为 0.2, 模拟的样本量为  $n = (100, 500, 1000)$ 。研究 4 种多重并行中介分析方法分别在上述 63 种模拟情形下的经验偏倚, 每种模拟情形重复 1000 次, 计算直接效应、特定中介效应的平均估计值、平均相对偏倚以及标准误; 并通过 1000 次 Bootstrap 计算直接效应和特定中介效应的 95% Bootstrap 置信区间, 从而计算直接效应和特定中介效应的估计精度, 以及 Bootstrap 法对 4 种多重并行中介分析方法得到的直接效应、特定中介效应发现非零效应的检验效能和覆盖率。

3. 模拟研究结果:模拟实验各样本量下结论基本一致, 以下介绍样本量为 1000 时的结果。

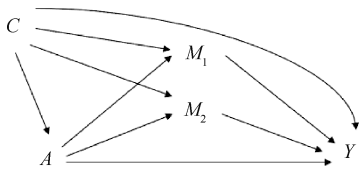


图 1 模拟实验的多重并行中介模型

(1)模拟情境一:当两个中介变量为连续型变量时,4种方法中纯回归法估计直接效应的平均相对偏倚最小,基于权重的填补法估计的平均相对偏倚最大。而当两个中介变量均为二分类变量时,4种方法中逆概率加权法估计直接效应的平均相对偏倚最小,基于权重的填补法估计的平均相对偏倚最大。见表4,5。

在估计间接效应时,不论中介变量为连续型变量或是二分类变量,4种方法中均是纯回归法对间接效应估计的偏倚最小,基于权重的填补法估计的

偏倚最大。

虽然基于权重的填补法在连续型中介情形下估计中介效应的偏倚小于二分类中介的情形,但是估计的偏倚均较大(>110%);95%CI覆盖真实值的概率为0,检验效能为100%。

当中介为二分类变量时,扩展的自然效应模型估计直接效应的偏倚小于中介为连续型变量时的估计偏倚,但对两个间接效应的估计偏倚明显大于中介为连续型变量时的估计。

(2)模拟情境二:在不同的中介相关性下,基于权重的填补法与扩展的自然效应模型相比于纯回归法和逆概率加权法,对各个效应的估计偏倚更大。且基于权重的填补法和扩展的自然效应模型的偏倚和系数有关。

(3)模拟情境三:在两个中介效应差不同的情形下,相比于纯回归法和逆概率加权法而言,基于

表 3 情境参数设置

情境	中介变量类型	$\rho(e_1, e_2)$	(IE1, IE2)
一	两个连续型变量或两个二分类变量	0	(0.4, 0.4)
二	两个连续型变量	-0.9, -0.6, -0.3, -0.1, 0, 0.1, 0.3, 0.6, 0.9	(0.4, 0.4)
三	两个连续型变量	0	(0.5, 0.7), (0.5, 1.0), (0.5, 1.3), (0.5, 1.8), (0.7, 1.0), (0.7, 1.3), (0.7, 1.8), (1.0, 1.3), (1.0, 1.8), (1.3, 1.8)

表 4 情境一中4种方法对直接效应的估计结果(n=1 000)

类别	方法	直接效应					
		平均估计值	平均相对偏倚(%)	标准误	估计精度	覆盖率(%)	检验效能(%)
连续型	纯回归法	0.200	0.0	0.044	0.091	93.4	94.7
	逆概率加权法	0.204	2.0	0.049	0.101	96.0	94.7
	基于权重的填补法	0.779	289.5	0.025	0.050	0.0	100.0
	扩展的自然效应模型	0.314	57.0	0.069	0.145	57.4	98.7
二分类	纯回归法	0.198	-1.0	0.034	0.067	95.1	95.4
	逆概率加权法	0.200	0.0	0.034	0.069	95.6	95.6
	基于权重的填补法	0.715	257.5	0.022	0.044	0.0	100.0
	扩展的自然效应模型	0.227	13.5	0.038	0.076	89.7	94.1

表 5 情境一中4种方法对间接效应的估计结果(n=1 000)

类别	方法	第一个中介效应						第二个中介效应					
		平均估计值	平均相对偏倚(%)	标准误	估计精度	覆盖率(%)	检验效能(%)	平均估计值	平均相对偏倚(%)	标准误	估计精度	覆盖率(%)	检验效能(%)
连续型	纯回归法	0.400	0.0	0.044	0.090	96.5	95.5	0.394	-1.5	0.044	0.090	96.2	95.1
	逆概率加权法	0.389	-2.8	0.059	0.116	94.8	94.5	0.376	-6.0	0.059	0.115	92.9	95.6
	基于权重的填补法	0.874	118.5	0.014	0.028	0.0	100.0	0.873	118.3	0.016	0.030	0.0	100.0
	扩展的自然效应模型	0.513	28.3	0.045	0.088	23.1	100.0	0.503	25.8	0.046	0.092	32.5	99.6
二分类	纯回归法	0.409	2.3	0.091	0.189	94.0	95.9	0.406	1.5	0.079	0.164	94.5	94.0
	逆概率加权法	0.411	2.8	0.094	0.199	95.0	95.4	0.411	2.8	0.083	0.173	94.7	95.3
	基于权重的填补法	0.968	142.0	0.008	0.015	0.0	100.0	0.959	139.8	0.009	0.017	0.0	100.0
	扩展的自然效应模型	0.842	110.5	0.038	0.077	0.0	100.0	0.810	102.5	0.038	0.077	0.0	100.0

权重的填补法和扩展的自然效应模型对各效应的估计偏倚均较大。但是在估计中介效应时,当不存在中介效应( $IE=1$ )时,基于权重的填补法和扩展的自然效应模型对其估计偏倚较小。见图 2。

### 三、实例应用

UK Biobank 是一项前瞻性队列研究,旨在调查

一系列疾病的遗传因素、生活方式和环境因素<sup>[20-22]</sup>。该研究共招募了 502 656 名年龄在 40~69 岁的受试者<sup>[23-24]</sup>。

既往研究表明,较高的水果摄入量与较低的心血管疾病风险有关<sup>[25-27]</sup>。而水果的摄入很大程度上也会影响血糖、血脂等血生化指标<sup>[25,28-29]</sup>。因此,

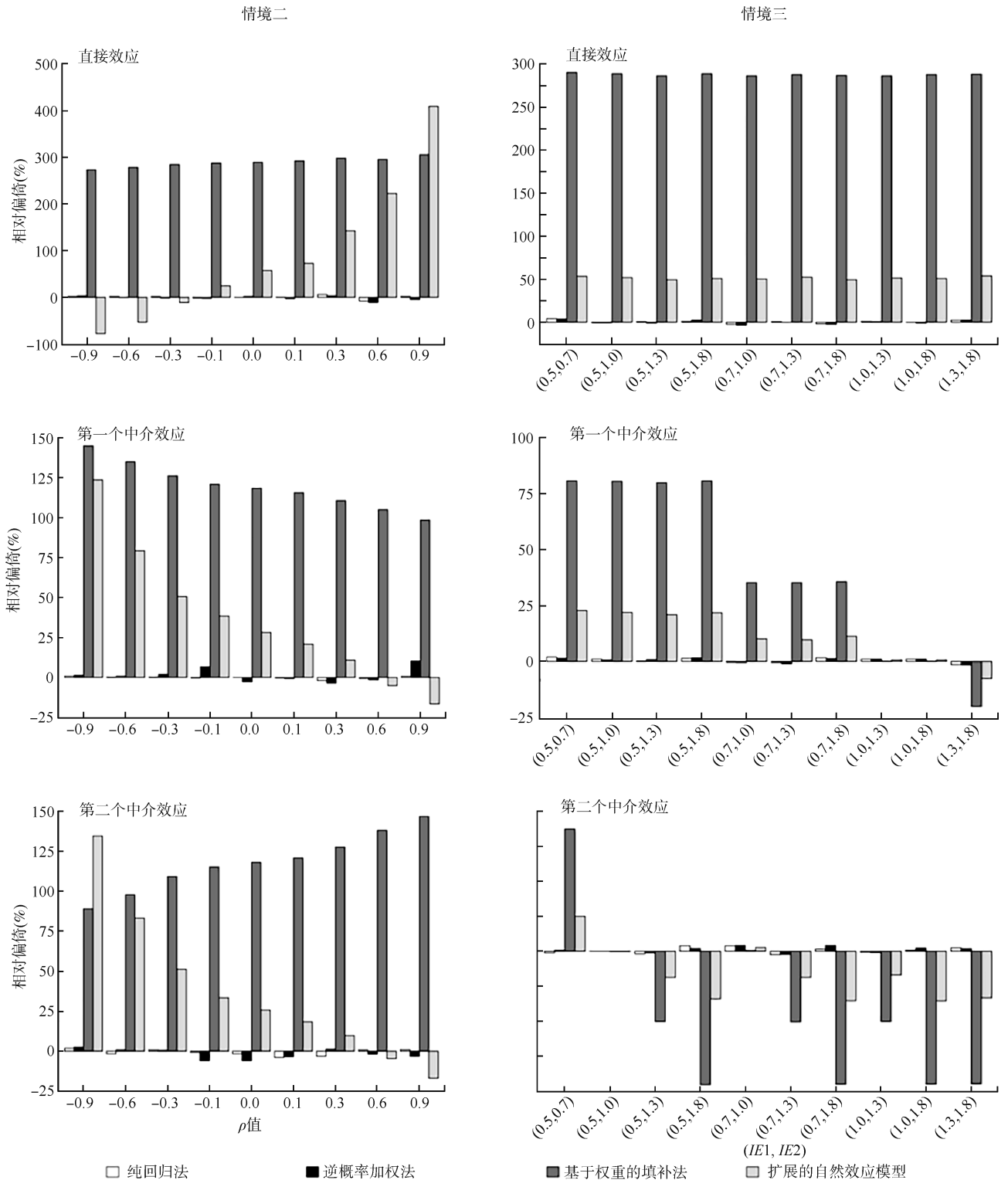


图 2 4 种方法对各效应的估计偏倚( $n=1\ 000$ )

本研究利用 UK Biobank 数据,除了研究水果摄入对心血管疾病风险的直接效应外,还探索了两个中介变量血糖和胆固醇是否在其因果路径中起作用以及中介效应的大小。本实例中暴露为新鲜水果摄入(二分类变量),结局为心血管疾病患病(二分类变量),两个并行中介变量血糖和胆固醇均为连续型变量。采用上述 4 种多重并行中介分析方法分别估计水果摄入对心血管疾病的直接效应、通过影响血糖水平而产生的间接效应以及通过影响胆固醇水平而产生的间接效应,同时调整年龄、性别、文化程度、经济水平、区域、饮酒、吸烟、运动和 BMI 等人口学特征变量。所有分析采用 R 3.6.3 软件进行。

除基于权重的填补法外,纯回归法、逆概率加权法以及扩展的自然效应模型估计的结果近似,即水果摄入与心血管疾病的风险降低有关,直接效应的大小分别为 0.685 (95%CI: 0.569~0.823)、0.654 (95%CI: 0.531~0.805)、0.677 (95%CI: 0.571~0.803)。此外,水果摄入可能会通过提高胆固醇水平而增加患心血管疾病的风险,纯回归法、逆概率加权法和扩展的自然效应模型估计的胆固醇间接效应分别为 1.017 (95%CI: 1.009~1.025)、1.020 (95%CI: 1.005~1.035) 和 1.016 (95%CI: 1.009~1.024),但这 3 种方法估计的血糖间接效应没有统计学意义。而基于权重的填补法相比于另外 3 种方法的结果相差较大,不仅识别不出水果摄入的直接效应 ( $OR=1.000$ ) 和胆固醇的间接效应 ( $OR=1.000$ ),还表明水果摄入可能会通过影响血糖水平而降低心血管疾病的风险 ( $OR=0.998$ , 95%CI: 0.997~0.999)。见表 6。

#### 四、讨论

多重中介分析研究的是自变量与因变量之间存在多个中介变量的情形。本文介绍了 4 种多重并行中介分析方法,即纯回归法、逆概率加权法、扩展的自然效应模型和基于权重的填补法,针对并行中介的情形,通过模拟实验比较了 4 种方法在不同的情境下,估计直接效应和间接效应的表现。

模拟结果表明,纯回归法在两个中介变量为连

续型变量时的表现优于两个二分类中介的情形,而逆概率加权法在两个中介变量为二分类变量时的表现优于两个连续型中介的情形,因此纯回归法推荐用于中介变量为连续型变量的情形,逆概率加权法推荐用于中介变量为二分类变量的情形。扩展的自然效应模型在两个二分类中介时估计直接效应的偏倚小于间接效应的估计偏倚,而在两个连续型中介时结果相反。而基于权重的填补法在不同的中介变量类型下对效应的估计偏倚均较大。

在两个并行中介均为连续型变量的设定下,两个中介变量间残差的相关性对纯回归法、逆概率加权法估计各效应的影响都不大,即纯回归法和逆概率加权法对各效应的估计偏倚都较小。基于权重的填补法对各效应的估计偏倚都较大。但当两个中介变量间的残差呈正相关时,扩展的自然效应模型估计间接效应的偏倚较小,当两个中介变量间残差的相关程度较小时,扩展的自然效应模型估计直接效应的偏倚较小。因而扩展的自然效应模型在用于两个并行中介的残差呈正相关且相关程度较小时更佳。

在不同的中介效应设置下,当中介效应 $<1$ 时,纯回归法和逆概率加权法的估计偏倚较小。基于权重的填补法仅在中介效应不存在( $IE1/2=1$ )时估计偏倚较小,而在其他情形下对各效应的估计偏倚相比于其他 3 种方法均较大。这可能是由于该方法最初是针对顺序中介的情形提出的<sup>[16]</sup>,而本文将用于两个并行中介的情形,因而可能存在较大的估计偏倚。提示在进行多重中介分析时,基于权重的填补法并不适用于并行中介的情形,研究者应当慎重选择分析方法。

虽然扩展的自然效应模型是用于多个并行中介的情形,但它要求中介变量间相互独立<sup>[13]</sup>。本研究的模拟实验中,也观察到扩展的自然效应模型在两个中介变量残差的相关程度较弱时表现相对较好。这提示研究者在进行多重并行中介分析时应当认真考虑中介变量间相关性的影响。

当中介变量为连续型变量时,可选用纯回归法;当中介变量为二分类变量时,可选用逆概率加

表 6 胆固醇和血糖在水果摄入与心血管疾病关联间的中介效应分析[OR 值(95%CI)]

方法	直接效应	胆固醇的间接效应	血糖的间接效应
纯回归法	0.685(0.569~0.823)	1.017(1.009~1.025)	1.001(0.998~1.003)
逆概率加权法	0.654(0.531~0.805)	1.020(1.005~1.035)	1.000(0.998~1.003)
基于权重的填补法	1.000(1.000~1.000)	1.000(1.000~1.000)	0.998(0.997~0.999)
扩展的自然效应模型	0.677(0.571~0.803)	1.016(1.009~1.024)	1.002(0.993~1.010)



权法。在中介间残差不同相关性的情形下,纯回归法和逆概率加权法表现不相上下,两种方法均可选用。由于扩展的自然效应模型直接参数化了直接效应和间接效应,因而若是出于效应估计的简便性,在两个并行中介的残差呈正相关且相关程度较小时可以考虑该方法。在并行中介的情形下,只有当中介效应不存在( $IE1/2=1$ )时,选用扩展的自然效应模型是优于纯回归法和逆概率加权法的,其他中介效应的情形依旧是选用纯回归法和逆概率加权法更佳。因而,可以针对同一实际数据采用几种不同的分析方法以考察结果的稳健性。基于权重的填补法可能更适用于顺序中介变量的情形<sup>[16]</sup>。

本文存在局限性。首先,本研究的模拟情形不够全面。例如,在考虑不同的中介变量类型时,两个中介变量分别为连续型变量和二分类变量的组合情形尚未考虑;在考虑两个中介变量残差间不同的相关性时,中介变量为二分类时的相关性也尚未考虑。其次,为了模型和分析的简便性,在整个研究中没有考虑交互作用,包括中介间的交互以及暴露与中介间的交互。在后续研究中,一方面笔者将考虑弥补上述不足,另一方面将考虑可以将本研究拓展到更多中介变量以及存在顺序中介变量的情形。

利益冲突 所有作者声明无利益冲突

作者贡献声明 于洋:统计分析、论文撰写;仇沁晓、尤东方、赵杨:研究指导、论文修改、经费支持

### 参 考 文 献

- Baron RM, Kenny DA. The moderator-mediator variable distinction in social psychological research: conceptual, strategic, and statistical considerations[J]. *J Pers Soc Psychol*, 1986, 51(6): 1173-1182. DOI: 10.1037/0022-3514.51.6.1173.
- Robins JM, Greenland S. Identifiability and exchangeability for direct and indirect effects[J]. *Epidemiology*, 1992, 3(2): 143-155. DOI:10.1097/00001648-199203000-00013.
- Imai K, Keele L, Tingley D. A general approach to causal mediation analysis[J]. *Psychol Methods*, 2010, 15(4): 309-334. DOI:10.1037/a0020761.
- Vanderweele TJ, Vansteelandt S. Odds ratios for mediation analysis for a dichotomous outcome[J]. *Am J Epidemiol*, 2010, 172(12):1339-1348. DOI:10.1093/aje/kwq332.
- Jérolon A, Baglietto L, Birmelé E, et al. Causal mediation analysis in presence of multiple mediators uncausally related[J]. *Int J Biostat*, 2020, 17(2): 191-221. DOI: 10.1515/ijb-2019-0088.
- Lai EY, Shih S, Huang YT, et al. A mediation analysis for a nonrare dichotomous outcome with sequentially ordered multiple mediators[J]. *Stat Med*, 2020, 39(10):1415-1428. DOI:10.1002/sim.8485.
- Rayan-Gharra N, Balicer RD, Tadmor B, et al. Association between cultural factors and readmissions: the mediating effect of hospital discharge practices and care-transition preparedness[J]. *BMJ Qual Saf*, 2019, 28(11): 866-874. DOI:10.1136/bmjqs-2019-009317.
- Tun HM, Bridgman SL, Chari R, et al. Roles of birth mode and infant gut Microbiota in intergenerational transmission of overweight and obesity from mother to offspring[J]. *JAMA Pediatr*, 2018, 172(4): 368-377. DOI: 10.1001/jamapediatrics.2017.5535.
- Hayes AF. Beyond baron and Kenny: statistical mediation analysis in the new millennium[J]. *Commun Monogra*, 2009, 76(4):408-420. DOI:10.1080/03637750903310360.
- MacKinnon DP. Introduction to statistical mediation analysis[M]. New York: Lawrence Erlbaum Associates, 2008.
- Preacher KJ, Hayes AF. Asymptotic and resampling strategies for assessing and comparing indirect effects in multiple mediator models[J]. *Behav Res Methods*, 2008, 40(3):879-891. DOI:10.3758/BRM.40.3.879.
- Vanderweele T, Vansteelandt S. Mediation analysis with multiple mediators[J]. *Epidemiol Methods*, 2014, 2(1): 95-115. DOI:10.1515/em-2012-0010.
- Lange T, Rasmussen M, Thygesen LC. Assessing natural direct and indirect effects through multiple pathways[J]. *Am J Epidemiol*, 2014, 179(4):513-518. DOI:10.1093/aje/kwt270.
- Taguri M, Featherstone J, Cheng J. Causal mediation analysis with multiple causally non-ordered mediators[J]. *Stat Methods Med Res*, 2018, 27(1): 3-19. DOI: 10.1177/0962280215615899.
- Lange T, Vansteelandt S, Bekaert M. A simple unified approach for estimating natural direct and indirect effects [J]. *Am J Epidemiol*, 2012, 176(3):190-195. DOI:10.1093/aje/kwr525.
- Steen J, Loeys T, Moerkerke B, et al. Flexible mediation analysis with multiple mediators[J]. *Am J Epidemiol*, 2017, 186(2):184-193. DOI:10.1093/aje/kwx051.
- Vansteelandt S, Bekaert M, Lange T. Imputation strategies for the estimation of natural direct and indirect effects[J]. *Epidemiol Methods*, 2012, 1(1): 131-158. DOI: 10.1515/2161-962X.1014.
- Wang W, Nelson S, Albert JM. Estimation of causal mediation effects for a dichotomous outcome in multiple-mediator models using the mediation formula [J]. *Stat Med*, 2013, 32(24): 4211-4228. DOI: 10.1002/sim.5830.
- Nguyen TQ, Webb-Vargas Y, Koning IM, et al. Causal mediation analysis with a binary outcome and multiple continuous or ordinal mediators: simulations and application to an alcohol intervention[J]. *Struct Equ Modeling*, 2016, 23(3):368-383. DOI:10.1080/10705511.2015.1062730.
- Allen N, Sudlow C, Downey P, et al. UK Biobank: current status and what it means for epidemiology[J]. *Health Policy Technol*, 2012, 1(3): 123-126. DOI: 10.1016/j.hlpt.2012.07.003.
- Manolio TA, Bailey-Wilson JE, Collins FS. Genes, environment and the value of prospective cohort studies [J]. *Nat Rev Genet*, 2006, 7(10): 812-820. DOI: 10.1038/nrg1919.
- Palmer LJ. UK Biobank: bank on it[J]. *Lancet*, 2007, 369(9578): 1980-1982. DOI: 10.1016/S0140-6736(07)60924.
- Allen NE, Sudlow C, Peakman T, et al. UK biobank data: come and get it[J]. *Sci Transl Med*, 2014, 6(224): 224ed4. DOI:10.1126/scitranslmed.3008601.
- Sudlow C, Gallacher J, Allen N, et al. UK biobank: an open access resource for identifying the causes of a wide range of complex diseases of middle and old age[J]. *PLoS Med*, 2015, 12(3): e1001779. DOI: 10.1371/journal.pmed.1001779.
- Du HD, Li LM, Bennett D, et al. Fresh fruit consumption and major cardiovascular disease in China[J]. *N Engl J Med*, 2016, 374(14): 1332-1343. DOI: 10.1056/NEJMoa1501451.
- Wang J, Liu FC, Li JX, et al. Fruit and vegetable consumption, cardiovascular disease, and all-cause mortality in China[J]. *Sci China Life Sci*, 2022, 65(1): 119-128. DOI:10.1007/s11427-020-1896-x.
- Zurbau A, Au-Yeung F, Mejia SB, et al. Relation of different fruit and vegetable sources with incident cardiovascular outcomes: a systematic review and meta-analysis of prospective cohort studies[J]. *J Am Heart Assoc*, 2020, 9(19):e017728. DOI:10.1161/JAHA.120.017728.
- Bragg F, Li LM, Bennett D, et al. Association of random plasma glucose levels with the risk for cardiovascular disease among Chinese adults without known diabetes[J]. *JAMA Cardiol*, 2016, 1(7): 813-823. DOI: 10.1001/jamacardio.2016.1702.
- Jiang ZL, Sun TY, He Y, et al. Dietary fruit and vegetable intake, gut microbiota, and type 2 diabetes: results from two large human cohort studies[J]. *BMC Med*, 2020, 18(1): 371. DOI:10.1186/s12916-020-01842-0.