

Choosing Wisely: Evaluating Latent Factor Models in the Presence of Contaminated Instrumental Variable(s) with Differing Strengths

Shubham Das¹ and Souvik Banerjee²

¹*PhD Candidate, Department of Economics, Indian Institute of Technology Bombay*

²*Assistant Professor, Department of Economics, Indian Institute of Technology Bombay*

Abstract

Background and Motivation: Causal inference methods are widely used in empirical research; however, there is a paucity of evidence on the use of latent factor methods in the presence of a contaminated instrumental variable (IV) for differing IV strength. We aim to build on and extend the work of Banerjee and Basu (2021) in multiple ways: (i) we discuss the identification criteria for the shared latent factor (SLF) and SLF+IV models, (ii) we provide a theoretical framework to illustrate how IV strength and contamination together determine the optimal choice between Two Stage Least Squares (2SLS), SLF, and SLF+IV estimators, (iii) we compare the finite sample properties of the naïve Ordinary Least Squares (OLS), SLF, 2SLS, and SLF+IV estimators in terms of percentage bias, standard error, coverage probability, and Mean Squared Error (MSE) for different degrees of contamination and strength of the IV, (iv) we illustrate the applicability of the OLS, 2SLS, SLF, and SLF+IV estimators using survey data to assess the causal impact of obesity on different indicators of health status, and (v) we demonstrate the virtue of latent factor models to control for effects of omitted variable(s).

Methods: We present a theoretical formulation to depict how the strength of an IV and contamination simultaneously determine the optimal choice between traditional IV estimators and shared latent factor estimators. We performed Monte Carlo simulations with four outcome variables and an endogenous treatment variable, with sample sizes varying between 500 and 1000, and for 1000 iterations, to compare the finite sample properties of the OLS, 2SLS, SLF, and SLF+IV estimators. Finally, we exhibit the applicability of the proposed estimators to study the causal impact of obesity on different health indicators (diastolic blood pressure reading, systolic blood pressure reading, blood glucose level, and haemoglobin level), using the arm circumference as an IV, and the 2019-2021 Round 5 of the National Family Health Survey (NFHS-5) data.

Results: We find that, for a given strength of the IV, there exists a threshold degree of contamination, such that the SLF+IV estimator has a lower (greater) bias than the SLF estimator when the degree of contamination lies below (above) that threshold. Similarly, we find that, for a given degree of contamination of the IV, there exists a threshold strength of the IV, such that the SLF+IV has a lower (greater) bias than the SLF estimator if the strength of the IV lies above (below) that threshold. The simulation results are consistent and align with our theoretical propositions. In our empirical application, we find that obesity is significantly associated with higher diastolic and systolic readings, a higher blood sugar level, and a higher haemoglobin level in the blood.

Conclusion: The choice between SLF and SLF+IV estimators based on the estimator bias is determined jointly by the strength of the IV and the degree of contamination. SLF+IV (SLF) is the optimal choice when the strength of the IV and the degree of contamination lie above (below) and below (above) certain thresholds, respectively.

Keywords: Causal inference, Latent factor models, Instrumental variable, Instrument strength, Instrument contamination, Parameter identification

JEL Codes: C3, C31, J13

1 Introduction

Causal inference is embodied in many empirical research works in economics. A researcher's objective is to obtain the treatment effect estimate which is much closer to the true treatment effect. Ordinary Least Squares (OLS) is a traditional method for estimating treatment effects. However, there can be a potential *endogeneity* in an OLS model. Endogeneity occurs due to: (i) the presence of unobserved factors affecting both the treatment and the outcome, (ii) measurement error, and (iii) reverse causality. The Two Stage Least Square (2SLS) is a traditional method wherein a special variable called the Instrumental Variable (IV) is employed that generates exogenous variation in the treatment variable, in order to mitigate the endogeneity problem. However, even if the IV is *strong*, i.e. it is highly correlated with the treatment variable, there can be a potential unobserved factor(s) which is(are) correlated with both the IV and the outcome, thus making the IV *contaminated*. If the degree of contamination of the IV is very high, then the 2SLS is likely to be more biased than the OLS estimator (Crown *et al.*, 2011; Nelson and Startz, 1990). Bound *et al.* (1993) calls this situation: *a cure, i.e. 2SLS analysis using a highly contaminated IV, worse than the disease, i.e. OLS with endogeneity*.

Focusing on the endogeneity problem arising due to the presence of unobservable factors affecting both the treatment and outcome variables, which is a major challenge to causal inference in much empirical research, latent factor models can be employed as an alternative to the traditional models. Latent factor models employ an exogenous latent factor that controls for such unobserved factors. These models can also be further extended by introducing an additional IV along with the exogenous latent factor. These models draw motivation from *path analysis*- a method that portrays how different observable and latent variables are related to each other (Goldberger, 1972; Hauser and Goldberger, 1971; Land, 1969). The idea of path analysis was first propounded by Wright (1921). The path analysis consisting of latent variables can be bifurcated into two branches: reflective modelling and formative modelling. In reflective modelling, an exogenous latent factor determines all the observed outcomes, while in formative modelling, all the observed outcomes determine the endogenous latent factor (Posey *et al.*, 2014; Lee *et al.*, 2013; Diamantopoulos and Siguaw, 2006; Diamantopoulos and Winklhofer, 2001; Law and Wong, 1999; Bollen and Lennox, 1991).

Estimation of latent factor models relies on the technique of Structural Equation Model (SEM) (Danes and Mann, 1984; Anderson and Gerbing, 1982). A major challenge in the application of SEM is *parameter identification*. The SEM is said to be identified if each unobserved parameter can be written as a function of the covariances between observed variables (Stein *et al.*, 2011; Lei and Wu, 2007; Mueller, 1997). It is often difficult to establish relationships between different variables in SEM in the case of underidentification (Bollen and Davis, 2009). There is no fixed rule for SEM identification, whereas the identification criteria depend on how the SEM has been specified (Bollen and Davis, 2009; Bentler and Weeks, 1980).

Banerjee and Basu (2021), published in *Econometrics* Vol(9), proposed two variants of latent factor models, viz. Shared Latent Factor without IV (SLF) and the Shared Latent Factor with IV (SLF+IV). In the SLF model, there is a group of outcome variables which are affected by a common treatment variable, and an exogenous latent factor is employed, which affects the treatment and all outcome variables. The SLF+IV model is an extension of the SLF model, wherein an additional IV is introduced that affects the treatment variable. Both models are variants of reflective modelling and are estimated using the technique of SEM. Banerjee and Basu highlighted the formulation of SLF and SLF+IV models, the minimum outcomes that are required for the identification of the SLF model, and the comparison among different estimators based on the estimator bias in the presence of a fairly strong IV with differing degrees of contamination. However, the paper leaves some questions unanswered, viz. (i) how to make a choice *between* SLF and SLF+IV estimators when both the IV strength and degree of contamination are differing, (ii) how the minimum number of outcomes

required for the identification of the SLF model has been ascertained, and (iii) number of outcomes that are required for the identification of the SLF+IV model.

In this paper, we build on and extend the work of Banerjee and Basu by (i) elucidating the analogy through which the minimum number of outcomes required for the identification of the SLF model has been ascertained, (ii) ascertaining the minimum number of outcomes required for the identification of the SLF+IV model, (ii) developed theoretical expositions to highlight how to make an optimal choice between SLF and SLF+IV when both IV strength and degree of contamination are differing, and (iv) comparing the finite sample properties of the OLS, SLF, 2SLS, and SLF+IV estimators in terms of percentage bias, standard error, coverage probability, and Mean Squared Error (MSE) when both IV strength and degree of contamination are differing, using Monte Carlo simulations. We further highlighted the applicability of the latent factor estimators, along with traditional estimators, through an empirical application, wherein we studied the causal impact of obesity on different indicators of health status using the nationally representative survey data from the latest round of the National Family Health Survey (NFHS-5) in India. Using the empirical application, we have illustrated how the latent factor in SLF and SLF+IV models captures the effects of omitted variable(s) by deliberately omitting a significant control variable from both models.

2 Materials and Methods

2.1 OLS Analysis

Let D be the treatment variable and Y be the outcome variable. The OLS model is described as follows:

$$Y = \beta_0 + \beta_1 D + v \quad (1)$$

where, β_0 is the constant term, β_1 captures the treatment effect, and v is the idiosyncratic error term. The bias of the OLS estimator is given as:

$$Bias(\hat{\beta}_{1,OLS}) = \frac{Cov(D, v)}{\sigma_D^2} \quad (2)$$

where, $Cov(D, v)$ is the covariance between D and v , and σ_D^2 is the variance of D . Ideally, D and v should be orthogonal, i.e. $Cov(D, v) = 0$ such that $\hat{\beta}_1$ is unbiased. However, if they are not orthogonal, then $\hat{\beta}_1$ may be biased.

2.2 SLF Analysis

Consider a system of m outcome variables Y^1, Y^2, \dots, Y^m , such that each outcome is a linear function of the treatment variable D and error terms v^1, v^2, \dots, v^m , respectively. D is the common treatment variable for all outcome variables. The effect of D on outcome variable $Y^i, i = 1, 2, \dots, m$ is captured by the parameter $\beta_1^i, i = 1, 2, \dots, m$, respectively.

$$Y^1 = \beta_1^1 D + v^1 \quad (3a)$$

$$Y^2 = \beta_1^2 D + v^2 \quad (3b)$$

⋮

$$Y^m = \beta_1^m D + v^m \quad (3m)$$

Let D be defined as follows:

$$D = \delta\theta + u \quad (4)$$

Here, θ is the latent factor. Now, assuming that the error terms of equations (3a), (3b), ..., (3m) are also linear functions of θ and an error term $w^i, i = 1, 2, \dots, m$, respectively. So, the composite error terms v^1, v^2, \dots, v^m are written as:

$$v^1 = \gamma^1\theta + w^1 \quad (5a)$$

$$v^2 = \gamma^2\theta + w^2 \quad (5b)$$

⋮

$$v^m = \gamma^m\theta + w^m \quad (5m)$$

Substituting equations (5a), (5b), ..., (5m) in equations (3a), (3b), ..., (3m), respectively, and writing all the equations together:

$$D = \delta\theta + u \quad (4)$$

$$Y^1 = \beta_1^1 D + \gamma^1\theta + w^1 \quad (6a)$$

$$Y^2 = \beta_1^2 D + \gamma^2\theta + w^2 \quad (6b)$$

⋮

$$Y^m = \beta_1^m D + \gamma^m\theta + w^m \quad (6m)$$

In the above system of equations, the latent factor, θ , is shared across the treatment and all the outcome equations, so θ is called the *shared latent factor*, and this model is called the *SLF* model. Since θ does not have any natural scale of measurement, we either need to set any one of the coefficients of θ to a constant and estimate $Var(\theta)$ or estimate all the coefficients of θ but set $Var(\theta)$ to a constant to identify the model. We set $\delta = 1$ without loss of generality. The following assumptions have been made regarding the shared latent factor and error terms:

$$Cov(\theta, u) = 0 \quad (7a)$$

$$Cov(\theta, w^h) = 0, h = 1, 2, \dots, m \quad (7b)$$

$$Cov(w^h, w^k) = 0, h \neq k, h = 1, 2, \dots, m, k = 1, 2, \dots, m \quad (7c)$$

$$Cov(u, w^h) = 0, h = 1, 2, \dots, m \quad (7d)$$

Assumption (7a) states that the shared latent factor is uncorrelated with the error term of the treatment variable equation. Assumption (7b) states that the shared latent factor is uncorrelated with the error term of the outcome variable equation after controlling for the shared latent factor. Assumption (7c) states that the error terms of the outcome variable equations, after controlling for the shared latent factor, are mutually uncorrelated. Finally, assumption (7d) states that the error term of the outcome variable equation, after controlling for the shared latent factor, is uncorrelated with the error term of the treatment variable equation. Assumptions (7a) to (7d) draw motivation from Diamantopoulos and Siguaw (2006) and Bollen and Lennox (1991).

Now, covariance equations have been computed, taking covariances between discordant pairs of outcome variables and covariances between each outcome variable and the treatment variable, considering assumptions (7a), (7b), (7c), and (7d). The covariance equations are as follows:

$$\begin{aligned}
Cov(Y^1, Y^2) &= Cov(\beta_1^1 D + \gamma^1 \theta + w^1, \beta_1^2 D + \gamma^2 \theta + w^2) \\
&= \beta_1^1 \beta_1^2 \sigma_D^2 + [\beta_1^1 \gamma^2 + \beta_1^2 \gamma^1 + \gamma^1 \gamma^2] \sigma_\theta^2 \\
Cov(Y^1, Y^3) &= Cov(\beta_1^1 D + \gamma^1 \theta + w^1, \beta_1^3 D + \gamma^3 \theta + w^3) \\
&= \beta_1^1 \beta_1^3 \sigma_D^2 + [\beta_1^1 \gamma^3 + \beta_1^3 \gamma^1 + \gamma^1 \gamma^3] \sigma_\theta^2 \\
&\vdots \\
Cov(Y^{m-1}, Y^m) &= Cov(\beta_1^{m-1} D + \gamma^{m-1} \theta + w^{m-1}, \beta_1^m D + \gamma^m \theta + w^m) \\
&= \beta_1^{m-1} \beta_1^m \sigma_D^2 + [\beta_1^{m-1} \gamma^m + \beta_1^m \gamma^{m-1} + \gamma^{m-1} \gamma^m] \sigma_\theta^2 \\
Cov(D, Y^1) &= Cov(D, \beta_1^1 D + \gamma^1 \theta + w^1) \\
&= \beta_1^1 \sigma_D^2 + \gamma^1 \sigma_\theta^2 \\
Cov(D, Y^2) &= Cov(D, \beta_1^2 D + \gamma^2 \theta + w^2) \\
&= \beta_1^2 \sigma_D^2 + \gamma^2 \sigma_\theta^2 \\
&\vdots \\
Cov(D, Y^m) &= Cov(\theta + u, \beta_1^m D + \gamma^m \theta + w^m) \\
&= \beta_1^m \sigma_D^2 + \gamma^m \sigma_\theta^2
\end{aligned}$$

where, σ_D^2 represents $Var(D)$ and σ_θ^2 represents $Var(\theta)$. In the above covariance equations, $Cov(Y^1, Y^2)$, $Cov(Y^1, Y^3)$, ..., $Cov(Y^{m-1}, Y^m)$, $Cov(D, Y^1)$, $Cov(D, Y^2)$, ..., $Cov(D, Y^m)$, and σ_D^2 are observable terms, while β_1^1 , β_1^2 , ..., β_1^m , γ^1 , γ^2 , ..., γ^m , and σ_θ^2 are parameters to be estimated. The general form of the covariance equations can be written as:

$$\begin{aligned}
Cov(Y^h, Y^k) &= \beta_1^h \beta_1^k \sigma_D^2 + [\beta_1^h \gamma^k + \beta_1^k \gamma^h + \gamma^h \gamma^k] \sigma_\theta^2, h \neq k, h = 1, 2, \dots, m, k = 1, 2, \dots, m \\
Cov(D, Y^h) &= \beta_1^h \sigma_D^2 + \gamma^h \sigma_\theta^2, h = 1, 2, \dots, m
\end{aligned}$$

The parameters of the SLF model can be identified only when a certain minimum number of outcome variables are available. There are m “ β ” parameters, m “ γ ” parameters, and 1 “ σ_θ^2 ” parameter, so there are a total of $m + m + 1 = 2m + 1$ parameters. There are $\frac{m(m-1)}{2}$ covariance terms between discordant pairs of outcomes and m covariance terms between the treatment and each outcome, so there are a total of $\frac{m(m-1)}{2} + m$ covariance terms. So, the minimum number of outcome variables required to identify all the model parameters can be determined by solving the following inequality (Carneiro *et al.*, 2003):

$$\begin{aligned}
\text{Number of covariance terms} &\geq \text{Number of parameters to be estimated} \\
\implies \frac{m(m-1)}{2} + m &\geq 2m + 1 \\
\implies m^2 - 3m - 2 &\geq 0
\end{aligned} \tag{8}$$

Solving inequality (8), we get $m \geq 4$ (approximately). So, a minimum of 4 outcome variables are needed to estimate all the structural parameters of the SLF model. When $m = 4$, there will be 10 covariance terms and 9 parameters. So, it will be an *overidentified* system of equations.

2.3 2SLS Analysis

In 2SLS analysis, the treatment variable is assumed to be endogenous such that it is a function of an IV. Let D be the treatment variable, \hat{D} be the predicted value of D , z be the IV, and Y be the outcome variable. The 2SLS model is described as follows:

$$D = \pi_0 + \pi_1 z + \tau \quad (9a)$$

$$Y = \beta_0 + \beta_1 \hat{D} + v \quad (9b)$$

where, π_0 and β_0 are constant terms, π_1 is the effect of z on D , β_1 captures the treatment effect, and τ and v are the idiosyncratic error terms. The bias of the 2SLS estimator is given as follows:

$$Bias(\hat{\beta}_{1,2SLS}) = \frac{Cov(z, v)}{Cov(D, z)} \quad (10)$$

where, $Cov(D, z)$ is the covariance between D and z (IV strength), and $Cov(z, v)$ is the covariance between z and v (degree of contamination of the IV). Ideally, the magnitude of $Cov(D, z)$ should be very high, and the value of $Cov(z, v)$ should be zero so that $\hat{\beta}_1$ may be unbiased. However, if the magnitude of $Cov(D, z)$ is very low (weak IV), and/or $Cov(z, v)$ is non-zero (contaminated IV), then $\hat{\beta}_1$ may be biased.

2.4 SLF+IV Analysis

The SLF+IV model is an extension of the SLF model, wherein an IV is introduced for the endogenous treatment variable. Let the IV again be denoted by z . The SLF+IV model is given as follows:

$$D = \pi z + \delta \theta + u \quad (11)$$

$$Y^1 = \beta_1^1 D + \gamma^1 \theta + w^1 \quad (12a)$$

$$Y^2 = \beta_1^2 D + \gamma^2 \theta + w^2 \quad (12b)$$

⋮

$$Y^m = \beta_1^m D + \gamma^m \theta + w^m \quad (12m)$$

In the above set of equations, the latent factor θ is again a shared latent factor as it is shared across the treatment variable and all the outcome variable equations. Since this model introduces an IV in the treatment variable equation of the SLF model, it is called the *SLF+IV model*. We again assume $\delta = 1$, without loss of generality. The SLF+IV model is also based on certain assumptions, which are as follows:

$$Cov(D, z) \neq 0 \quad (13a)$$

$$Cov(z, w^h) = 0, h = 1, 2, \dots, m \quad (13b)$$

$$Cov(z, u) = 0 \quad (13c)$$

$$Cov(\theta, u) = 0 \quad (13d)$$

$$Cov(\theta, w^h) = 0, h = 1, 2, \dots, m \quad (13e)$$

$$Cov(w^h, w^k) = 0, h \neq k, h = 1, 2, \dots, m, k = 1, 2, \dots, m \quad (13f)$$

$$Cov(u, w^h) = 0, h = 1, 2, \dots, m \quad (13g)$$

$$Cov(z, \theta) = 0 \quad (13h)$$

Assumption (13a) states that the treatment variable and the IV should be correlated i.e. the strength of the IV cannot be zero. Ideally, the correlation should be very high so that z is a strong IV. Assumption (13b) states that the IV should not be correlated with the error terms of the outcome variable equations after controlling for the shared latent factor, i.e. the IV should be uncontaminated. Assumption (13c) states that the IV should be uncorrelated with the error term of the treatment variable equation. Assumption (13d) states that the shared latent factor should be uncorrelated with the error term of the treatment variable equation. Assumption (13e) states that the shared latent factor should be uncorrelated with the error terms of the outcome variable equations after controlling for the shared latent factor. Assumption (13f) states that the error terms of the outcome variable equations should be mutually uncorrelated after controlling for the shared latent factor. Assumption (13g) states that the error term of the treatment variable equation should be uncorrelated with the error terms of the outcome variable equations after controlling for the shared latent factor. Finally, assumption (13h) states that the IV and the shared latent factor are uncorrelated. Again, assumptions (13a) to (13h) draw motivation from Diamantopoulos and Siguaw (2006) and Bollen and Lennox (1991).

Now, covariance equations have again been computed, taking covariances between discordant pairs of outcome variables, covariances between each outcome variable and the treatment variable, covariances between each outcome variable and the IV, and covariance between the treatment variable and the IV considering assumptions (13a), (13b), ..., (13h). The covariance equations are as follows:

$$\begin{aligned} Cov(Y^1, Y^2) &= Cov(\beta_1^1 D + \gamma^1 \theta + w^1, \beta_1^2 D + \gamma^2 \theta + w^2) \\ &= \beta_1^1 \beta_1^2 \sigma_D^2 + [\beta_1^1 \gamma^2 + \beta_1^2 \gamma^1 + \gamma^1 \gamma^2] \sigma_\theta^2 \\ Cov(Y^1, Y^3) &= Cov(\beta_1^1 D + \gamma^1 \theta + w^1, \beta_1^3 D + \gamma^3 \theta + w^3) \\ &= \beta_1^1 \beta_1^3 \sigma_D^2 + [\beta_1^1 \gamma^3 + \beta_1^3 \gamma^1 + \gamma^1 \gamma^3] \sigma_\theta^2 \\ &\vdots \\ Cov(Y^{m-1}, Y^m) &= Cov(\beta_1^{m-1} D + \gamma^{m-1} \theta + w^{m-1}, \beta_1^m D + \gamma^m \theta + w^m) \\ &= \beta_1^{m-1} \beta_1^m \sigma_D^2 + [\beta_1^{m-1} \gamma^m + \beta_1^m \gamma^{m-1} + \gamma^{m-1} \gamma^m] \sigma_\theta^2 \\ Cov(Y^1, D) &= Cov(\beta_1^1 D + \gamma^1 \theta + w^1, D) \\ &= \beta_1^1 \sigma_D^2 + \gamma^1 \sigma_\theta^2 \\ Cov(Y^2, D) &= Cov(\beta_1^2 D + \gamma^2 \theta + w^2, D) \\ &= \beta_1^2 \sigma_D^2 + \gamma^2 \sigma_\theta^2 \\ &\vdots \\ Cov(Y^m, D) &= Cov(\beta_1^m D + \gamma^m \theta + w^m, D) \\ &= \beta_1^m \sigma_D^2 + \gamma^m \sigma_\theta^2 \\ Cov(Y^1, z) &= Cov(\beta_1^1 D + \gamma^1 \theta + w^1, z) \\ &= \beta_1^1 \pi \sigma_z^2 \\ Cov(Y^2, z) &= Cov(\beta_1^2 D + \gamma^2 \theta + w^2, z) \\ &= \beta_1^2 \pi \sigma_z^2 \\ &\vdots \end{aligned}$$

$$\begin{aligned}
Cov(Y^m, z) &= Cov(\beta_1^m D + \gamma^m \theta + w^m, z) \\
&= \beta_1^m \pi \sigma_z^2 \\
Cov(D, z) &= \pi \sigma_z^2
\end{aligned}$$

In the above covariance equations, σ_D^2 represents $Var(D)$, σ_θ^2 represents $Var(\theta)$, and σ_z^2 represents $Var(z)$. $Cov(Y^1, Y^2)$, $Cov(Y^1, Y^3)$, ..., $Cov(Y^{m-1}, Y^m)$, $Cov(Y^1, D)$, $Cov(Y^2, D)$, ..., $Cov(Y^m, D)$, $Cov(Y^1, z)$, $Cov(Y^2, z)$, ..., $Cov(Y^m, z)$, $Cov(D, z)$, σ_z^2 , and σ_D^2 are observable terms, while $\beta_1^1, \beta_2^2, \dots, \beta_m^m, \gamma^1, \gamma^2, \dots, \gamma^m, \sigma_\theta^2$, and π_1 are parameters to be estimated. The general form of the covariance equations can be written as follows:

$$\begin{aligned}
Cov(Y^h, Y^k) &= \beta_1^h \beta_1^k \sigma_D^2 + [\beta_1^h \gamma^k + \beta_1^k \gamma^h + \gamma^h \gamma^k] \sigma_\theta^2, h \neq k, h = 1, 2, \dots, m, k = 1, 2, \dots, m \\
Cov(Y^h, D) &= \beta_1^h \sigma_D^2 + \gamma^h \sigma_\theta^2, h = 1, 2, \dots, m \\
Cov(Y^h, z) &= \beta_1^h \pi_1 \sigma_z^2, h = 1, 2, \dots, m \\
Cov(D, z) &= \pi_1 \sigma_z^2
\end{aligned}$$

Like the SLF model, all the parameters of the SLF+IV model can be identified only when a certain minimum number of outcome variables are available. There are m “ β ” parameters, m “ γ ” parameters, 1 “ π ” parameter, and 1 “ σ_θ^2 ” parameter, so there are a total of $m+m+1+1 = 2m+2$ parameters to be estimated. There are $\frac{m(m-1)}{2}$ covariance terms between discordant pairs of outcomes, m covariance terms between the treatment and each outcome variable, m covariance terms between the IV and each outcome variable, and 1 covariance term between the treatment variable and the IV, so there are a total of $\frac{m(m-1)}{2} + m + m + 1 = \frac{m(m-1)}{2} + 2m + 1$ covariance terms. So, the minimum number of outcome variables required to identify all the model parameters can be determined by solving the following inequality (Carneiro et al., 2003):

$$\begin{aligned}
\text{Number of covariance terms} &\geq \text{Number of parameters to be estimated} \\
\implies \frac{m(m-1)}{2} + 2m + 1 &\geq 2m + 2 \\
\implies m^2 - m - 2 &\geq 0
\end{aligned} \tag{14}$$

Solving the inequality (14) we get $m \geq 2$. So, the parameters of the SLF+IV model can be estimated when at least 2 outcome measurements are available¹, which are 2 outcomes fewer than that required to estimate the SLF model. When $m = 2$, there will be 6 covariance terms and 6 parameters. So, it will be a *just identified* system of equations.

2.5 Comparison Between Biases of SLF and SLF+IV

This section aims to elucidate how the degree of contamination and the strength of the IV determine the choice between SLF and SLF+IV estimators based on their absolute biases. In the following propositions, each variable has n observations.

Proposition 1: *The absolute bias of the SLF+IV estimator is lower (greater) than that of the SLF estimator when the degree of contamination of the IV lies below (above) a certain threshold for a given strength of the IV.*

¹ $m \geq 1$ when there are two IVs (See Appendix). So, taking two IVs reduces the number of outcome variables required for estimating the SLF+IV model.

Proof: Consider the following SLF model:

$$\begin{aligned}
D_j &= \alpha_{0j} + \theta_j + u_j \\
Y_j^1 &= \beta_{0j}^1 + \beta_{1j}^1 D_j + \gamma_j^1 \theta_j + w_j^1 \\
Y_j^2 &= \beta_{0j}^2 + \beta_{1j}^2 D_j + \gamma_j^2 \theta_j + w_j^2 \\
&\vdots \\
Y_j^m &= \beta_{0j}^m + \beta_{1j}^m D_j + \gamma_j^m \theta_j + w_j^m
\end{aligned}$$

where, $j = 1, 2, \dots, n$. In the above system of equations, treatment effect parameters are represented by $\beta_{1j}^1, \beta_{1j}^2, \dots, \beta_{1j}^m$. Let $\hat{\beta}_{1,SLF}^i$ be the estimated value of β_1^i , using some finite sample, where $\hat{\beta}_{1,SLF}^i$ is the treatment effect estimate for any arbitrary i^{th} outcome, such that $1 \leq i \leq m$. The expression for $\hat{\beta}_{1,SLF}^i$ is given as follows (Maddala 1992):

$$\begin{aligned}
\hat{\beta}_{1,SLF}^i &= \frac{[\sum_{j=1}^n (\theta_j - \bar{\theta})^2][\sum_{j=1}^n (D_j - \bar{D})(Y_j^i - \bar{Y}^i)] - [\sum_{j=1}^n (D_j - \bar{D})(\theta_j - \bar{\theta})][\sum_{j=1}^n (\theta_j - \bar{\theta})(Y_j^i - \bar{Y}^i)]}{[\sum_{j=1}^n (D_j - \bar{D})^2][\sum_{j=1}^n (\theta_j - \bar{\theta})^2] - [\sum_{j=1}^n (D_j - \bar{D})(\theta_j - \bar{\theta})]^2} \\
&= \frac{\sigma_\theta^2 \text{Cov}(D, Y^i) - \text{Cov}(D, \theta) \text{Cov}(\theta, Y^i)}{\sigma_D^2 \sigma_\theta^2 - [\text{Cov}(D, \theta)]^2} \\
&= \frac{\sigma_\theta^2 \text{Cov}(D, \beta_0^i + \beta_1^i D + \gamma^i \theta + w^i) - \text{Cov}(D, \theta) \text{Cov}(\theta, \beta_0^i + \beta_1^i D + \gamma^i \theta + w^i)}{\sigma_D^2 \sigma_\theta^2 - [\text{Cov}(\alpha_0 + \theta + u, \theta)]^2} \\
&= \frac{\sigma_\theta^2 [\beta_1^i \sigma_D^2 + \gamma^i \sigma_\theta^2 + \text{Cov}(D, w^i)] - \sigma_\theta^2 (\beta_1^i \sigma_\theta^2 + \gamma^i \sigma_\theta^2)}{\sigma_D^2 \sigma_\theta^2 - \sigma_\theta^4} \\
&= \frac{\beta_1^i \sigma_D^2 + \gamma^i \sigma_\theta^2 + \text{Cov}(D, w^i) - \beta_1^i \sigma_\theta^2 - \gamma^i \sigma_\theta^2}{\sigma_D^2 - \sigma_\theta^2} \\
&= \frac{\beta_1^i \sigma_D^2 + \text{Cov}(\alpha_0 + \theta + u, w^i) - \beta_1^i \sigma_\theta^2}{\sigma_D^2 - \sigma_\theta^2} \\
&= \frac{\beta_1^i (\sigma_D^2 - \sigma_\theta^2)}{\sigma_D^2 - \sigma_\theta^2} + \frac{\text{Cov}(u, w^i)}{\sigma_D^2 - \sigma_\theta^2} \\
&= \beta_1^i + \frac{\text{Cov}(u, w^i)}{\sigma_\theta^2 + \sigma_u^2 - \sigma_\theta^2} \\
&= \beta_1^i + \frac{\text{Cov}(u, w^i)}{\sigma_u^2} \tag{15}
\end{aligned}$$

So, equation (15) shows that the expression for the bias of the SLF estimator is given as:

$$\text{Bias}(\hat{\beta}_{1,SLF}^i) = \frac{\text{Cov}(u, w^i)}{\sigma_u^2} \tag{16}$$

We have relaxed the assumption that $\text{Cov}(u, w^i) = 0$, so $\text{Cov}(u, w^i)$ can take any values. Now, consider the following SLF+IV model:

$$\begin{aligned}
D_j &= \alpha_{0j} + \pi_j z_j + \theta_j + u_j \\
Y_j^1 &= \beta_{0j}^1 + \beta_{1j}^1 D_j + \gamma_j^1 \theta_j + w_j^1 \\
Y_j^2 &= \beta_{0j}^2 + \beta_{1j}^2 D_j + \gamma^2 \theta_j + w_j^2 \\
&\vdots \\
Y_j^m &= \beta_{0j}^m + \beta_{1j}^m D_j + \gamma_j^m \theta_j + w_j^m
\end{aligned}$$

where, $j = 1, 2, \dots, n$. In the above system of equations, treatment effect parameters are represented by $\beta_{1j}^1, \beta_{1j}^2, \dots, \beta_{1j}^m$. Let $\hat{\beta}_{1,SLF+IV}^i$ be the estimated value of β_1^i , using some finite sample, where the same i^{th} outcome has been considered that was taken for the SLF model, and $1 \leq i \leq m$. The expression for $\hat{\beta}_{1,SLF+IV}^i$ is given as follows (Maddala 1992):

$$\begin{aligned}
\hat{\beta}_{1,SLF+IV}^i &= \frac{[\sum_{j=1}^n (\theta_j - \bar{\theta})^2][\sum_{j=1}^n (D_j - \bar{D})(Y_j^i - \bar{Y}^i)] - [\sum_{j=1}^n (D_j - \bar{D})(\theta_j - \bar{\theta})][\sum_{j=1}^n (\theta_j - \bar{\theta})(Y_j^i - \bar{Y}^i)]}{[\sum_{j=1}^n (D_j - \bar{D})^2][\sum_{j=1}^n (\theta_j - \bar{\theta})^2] - [\sum_{j=1}^n (D_j - \bar{D})(\theta_j - \bar{\theta})]^2} \\
&= \frac{\sigma_\theta^2 \text{Cov}(D, Y^i) - \text{Cov}(D, \theta) \text{Cov}(\theta, Y^i)}{\sigma_D^2 \sigma_\theta^2 - [\text{Cov}(D, \theta)]^2} \\
&= \frac{\sigma_\theta^2 \text{Cov}(D, \beta_0^i + \beta_1^i D + \gamma^i \theta + w^i) - \text{Cov}(D, \theta) \text{Cov}(\theta, \beta_0^i + \beta_1^i D + \gamma^i \theta + w^i)}{\sigma_D^2 \sigma_\theta^2 - [\text{Cov}(\alpha_0 + \pi z + \theta + u, \theta)]^2} \\
&= \frac{\sigma_\theta^2 [\beta_1^i \sigma_D^2 + \gamma^i \pi \text{Cov}(z, \theta) + \gamma^i \sigma_\theta^2 + \text{Cov}(D, w^i)] - [\pi \text{Cov}(z, \theta) + \sigma_\theta^2][\beta_1^i \pi \text{Cov}(z, \theta) + \beta_1^i \sigma_\theta^2 + \gamma^i \sigma_\theta^2]}{\sigma_D^2 \sigma_\theta^2 - [\pi \text{Cov}(z, \theta) + \sigma_\theta^2]^2} \\
&= \frac{\sigma_\theta^2 \beta_1^i \sigma_D^2 + \sigma_\theta^2 \text{Cov}(D, w^i) - \beta_1^i \pi^2 [\text{Cov}(z, \theta)^2 - \beta_1^i \pi \sigma_\theta^2 - 2\beta_1^i \pi \sigma_\theta^2 \text{Cov}(z, \theta)]}{\sigma_D^2 \sigma_\theta^2 - [\pi \text{Cov}(z, \theta) + \sigma_\theta^2]^2} \\
&= \frac{\beta_1^i [\sigma_D^2 \sigma_\theta^2 - \{\pi \text{Cov}(z, \theta) + \sigma_\theta^2\}^2] + \sigma_\theta^2 \text{Cov}(D, w^i)}{\sigma_D^2 \sigma_\theta^2 - [\pi \text{Cov}(z, \theta) + \sigma_\theta^2]^2} \\
&= \beta_1^i + \frac{\sigma_\theta^2 \text{Cov}(u, w^i)}{\sigma_D^2 \sigma_\theta^2 - [\pi \text{Cov}(z, \theta) + \sigma_\theta^2]^2}
\end{aligned} \tag{17}$$

So, equation (17) shows that the expression for the bias of the SLF+IV estimator is given as:

$$\begin{aligned}
\text{Bias}(\hat{\beta}_{1,SLF+IV}^i) &= \frac{\sigma_\theta^2 \text{Cov}(u, w^i)}{\sigma_D^2 \sigma_\theta^2 - [\pi \text{Cov}(z, \theta) + \sigma_\theta^2]^2} \\
\implies \text{Bias}(\hat{\beta}_{1,SLF+IV}^i) &= \frac{\sigma_\theta^2 \text{Cov}(u, w^i)}{\sigma_\theta^2 [\pi^2 \sigma_z^2 + \sigma_\theta^2 + \sigma_u^2 + 2\pi \text{Cov}(z, \theta)] - [\pi \text{Cov}(z, \theta) + \sigma_\theta^2]^2} \\
\implies \text{Bias}(\hat{\beta}_{1,SLF+IV}^i) &= \frac{\sigma_\theta^2 \text{Cov}(u, w^i)}{\pi^2 \sigma_z^2 \sigma_\theta^2 + \sigma_\theta^2 \sigma_u^2 - \pi^2 [\text{Cov}(z, \theta)]^2}
\end{aligned} \tag{18}$$

Note that the term $|\text{Cov}(z, \theta)|$ represents the degree of contamination of the IV (Banerjee and Basu, 2021). We have relaxed the assumption $\text{Cov}(z, \theta) = 0$, so $\text{Cov}(z, \theta)$ can take any values apart from 0. We consider only the case wherein π and $\text{Cov}(z, \theta)$ have the same signs, i.e., either both positive or both are negative. In that case, the term σ_z represents the strength of the IV (see Appendix). In equation (18) if $\pi = 0$, then the bias of the SLF+IV estimator shall be the same as that of the SLF estimator. Now, using equation (16) in equation (18):

$$\begin{aligned} Bias(\hat{\beta}_{1,SLF+IV}^i) &= \frac{\sigma_\theta^2 \sigma_u^2 Bias(\hat{\beta}_{1,SLF}^i)}{\pi^2 \sigma_z^2 \sigma_\theta^2 + \sigma_\theta^2 \sigma_u^2 - \pi^2 [Cov(z, \theta)]^2} \\ \implies \frac{Bias(\hat{\beta}_{1,SLF+IV}^i)}{Bias(\hat{\beta}_{1,SLF}^i)} &= \frac{\sigma_\theta^2 \sigma_u^2}{\pi^2 \sigma_z^2 \sigma_\theta^2 + \sigma_\theta^2 \sigma_u^2 - \pi^2 [Cov(z, \theta)]^2} \end{aligned}$$

Taking modulus on both sides:

$$\begin{aligned} \left| \frac{Bias(\hat{\beta}_{1,SLF+IV}^i)}{Bias(\hat{\beta}_{1,SLF}^i)} \right| &= \left| \frac{\sigma_\theta^2 \sigma_u^2}{\pi^2 \sigma_z^2 \sigma_\theta^2 + \sigma_\theta^2 \sigma_u^2 - \pi^2 [Cov(z, \theta)]^2} \right| \\ \implies \frac{\left| Bias(\hat{\beta}_{1,SLF+IV}^i) \right|}{\left| Bias(\hat{\beta}_{1,SLF}^i) \right|} &= \frac{|\sigma_\theta^2 \sigma_u^2|}{|\pi^2 \sigma_z^2 \sigma_\theta^2 + \sigma_\theta^2 \sigma_u^2 - \pi^2 [Cov(z, \theta)]^2|} \end{aligned} \quad (19)$$

From equation (19), two cases emerge if the absolute bias of the SLF+IV estimator has to be lower than or equal to that of the SLF estimator:

Case 1:

$$\begin{aligned} \pi^2 \sigma_z^2 \sigma_\theta^2 - \pi^2 [Cov(z, \theta)]^2 &\geq 0 \\ \implies [Cov(z, \theta)]^2 &\leq \sigma_z^2 \sigma_\theta^2 \\ \implies |Cov(z, \theta)| &\leq [\sigma_z \sigma_\theta]^* \end{aligned} \quad (20)$$

In inequality (20), the term $[\sigma_z \sigma_\theta]^*$ represents the threshold degree of contamination of the IV, such that if the degree of contamination of the IV lies below (above) the threshold, then the SLF+IV estimator will have a lower (greater) bias than the SLF estimator for a given strength of the IV (σ_z).

Case 2:

$$\begin{aligned} \sigma_\theta^2 \sigma_u^2 &\geq -\pi^2 \sigma_z^2 \sigma_\theta^2 - \sigma_\theta^2 \sigma_u^2 + \pi^2 [Cov(z, \theta)]^2 \\ \implies \pi^2 [Cov(z, \theta)]^2 &\leq \pi^2 \sigma_z^2 \sigma_\theta^2 + 2\sigma_\theta^2 \sigma_u^2 \\ \implies |\pi| |Cov(z, \theta)| &\leq \sigma_\theta \sqrt{\pi^2 \sigma_z^2 + 2\sigma_u^2} \\ \implies |Cov(z, \theta)| &\leq \left[\frac{\sigma_\theta \sqrt{\pi^2 (\sigma_z)^2 + 2\sigma_u^2}}{|\pi|} \right]^* \end{aligned} \quad (21)$$

In inequality (21), the term $\left[\frac{\sigma_\theta \sqrt{\pi^2 (\sigma_z)^2 + 2\sigma_u^2}}{|\pi|} \right]^*$ represents the threshold degree of contamination of the IV, such that if the degree of contamination of the IV lies below (above) the threshold, then the SLF+IV estimator will have a lower (greater) bias than the SLF estimator for a given strength of the IV (σ_z).

Hence, both cases 1 and 2 indicate the existence of a threshold degree of contamination of the IV that determines the choice between SLF and SLF+IV estimators, for a given strength of the IV. \square

Proposition 2: *The absolute bias of the SLF+IV estimator is lower (greater) than that of the SLF estimator when the strength of the IV is higher (lower) than a certain threshold level for a given degree of contamination.*

Proof: Consider the same formulations of SLF and SLF+IV models that were taken for proving the proposition 1. The term $|Cov(z, \theta)|$ again represents the degree of contamination. We have made the same assumption about the signs of π and $Cov(z, \theta)$ as made in Proposition 1, such that the term σ_z represents the strength of the IV. Revisiting equation (19):

$$\frac{|Bias(\hat{\beta}_{1,SLF+IV}^i)|}{|Bias(\hat{\beta}_{1,SLF}^i)|} = \frac{|\sigma_\theta^2 \sigma_u^2|}{|\pi^2 \sigma_z^2 \sigma_\theta^2 + \sigma_\theta^2 \sigma_u^2 - \pi^2 [Cov(z, \theta)]^2|}$$

From equation (19), again two cases emerge if the absolute bias of the SLF+IV estimator has to be lower than or equal to that of the SLF estimator:

Case 1:

$$\begin{aligned} \pi^2 \sigma_z^2 \sigma_\theta^2 - \pi^2 [Cov(z, \theta)]^2 &\geq 0 \\ \implies \sigma_z^2 &\geq \frac{[Cov(z, \theta)]^2}{\sigma_\theta^2} \\ \implies \sigma_z &\geq \left[\frac{|Cov(z, \theta)|}{|\sigma_\theta|} \right]^* \end{aligned} \tag{22}$$

In inequality (22), the term $\left[\frac{|Cov(z, \theta)|}{|\sigma_\theta|} \right]^*$ represents the threshold strength of the IV, such that the SLF+IV estimator will have a lower (greater) bias than the SLF estimator when the strength is greater (lower) than the threshold for a given degree of contamination of the IV ($|Cov(z, \theta)|$).

Case 2:

$$\begin{aligned} \sigma_\theta^2 \sigma_u^2 &\geq -\pi^2 \sigma_z^2 \sigma_\theta^2 - \sigma_\theta^2 \sigma_u^2 + \pi^2 [Cov(z, \theta)]^2 \\ \implies \pi^2 \sigma_z^2 \sigma_\theta^2 &\geq \pi^2 [Cov(z, \theta)]^2 - 2\sigma_\theta^2 \sigma_u^2 \\ \implies \sigma_z^2 &\geq \frac{\pi^2 [|Cov(z, \theta)|]^2 - 2\sigma_\theta^2 \sigma_u^2}{\pi^2 \sigma_\theta^2} \quad (\text{Since, } [Cov(z, \theta)]^2 = [|Cov(z, \theta)|]^2) \\ \implies \sigma_z &\geq \left[\frac{\left| \sqrt{\pi^2 [|Cov(z, \theta)|]^2 - 2\sigma_\theta^2 \sigma_u^2} \right|}{|\pi| |\sigma_\theta|} \right]^* \end{aligned} \tag{23}$$

In inequality (23), the term $\left[\frac{\left| \sqrt{\pi^2 [|Cov(z, \theta)|]^2 - 2\sigma_\theta^2 \sigma_u^2} \right|}{|\pi| |\sigma_\theta|} \right]^*$ represents the threshold strength of the IV, provided that $\pi^2 [|Cov(z, \theta)|]^2 \geq 2\sigma_\theta^2 \sigma_u^2$, such that the SLF+IV estimator will have a lower (greater)

bias than the SLF estimator when the strength is greater (lower) than the threshold for a given degree of contamination of the IV ($|Cov(z, \theta)|$).

Hence, in both cases 1 and 2, it is evident that there exists a threshold strength of the IV that determines the choice between the SLF and SLF+IV estimators for a given degree of contamination of the IV.

□

So, propositions 1 and 2 together indicate the following result:

If the degree of contamination lies below (above) a certain threshold and the strength of the IV lies above (below) a certain threshold, then the SLF+IV shall have a lower (greater) absolute bias than the SLF estimator.

3 Simulations

In this section, we aim to analyze thresholds for the IV strength and the degree of contamination to choose between SLF and SLF+IV estimators using Monte Carlo simulations. We shall compare SLF, SLF+IV, OLS, and 2SLS estimates based on bias, standard error, coverage probability, and Mean Squared Error (MSE).

3.1 Baseline Cases: OLS and SLF

We have considered OLS and SLF as the baseline cases. The baseline case considers continuous outcome variable $Y^i, i = 1, 2, 3, 4$, a binary treatment variable (D), an observed continuous control variable (X), a latent factor variable (θ), and a continuous IV (z). Taking four outcome variables shall enable us to identify all the structural parameters of SLF and SLF+IV models. The baseline case has been constructed such that z and θ are uncorrelated, and the standard deviation of z is 1 (Basu and Chan, 2014; Banerjee and Basu, 2021). The baseline case has been constructed for $n = 500$ observations. The Data Generating Process (DGP) for the treatment variable is as follows:

$$D^* = \beta_0 + \beta_1 X + \beta_2 z + \beta_3 \theta + \zeta$$

$$D = \begin{cases} 1, & \text{if } D^* > 0 \\ 0, & \text{if } D^* \leq 0 \end{cases}$$

We assume that the true values (model parameters) are as follows: $\beta_0 = 1$, $\beta_1 = 2$, $\beta_2 = 3$, and $\beta_3 = 2$. X & ζ are normally and independently and identically distributed (IID) with mean 0 and variance 1. z and θ have been drawn from a bivariate normal distribution, such that $\begin{pmatrix} \theta \\ z \end{pmatrix} \sim \mathcal{N}\left(\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 1 & \sigma_{z\theta} \\ \sigma_{z\theta} & \sigma_z * \sigma_z \end{pmatrix}\right)$ where, $\sigma_{z\theta} = Cov(z, \theta)$ and represents the degree of contamination, and σ_z represents the strength of the IV. In the baseline case, $\sigma_{z\theta}$ is 0 and σ_z is 1.

The outcome variable $Y^i, i = 1, 2, 3, 4$ has been generated as follows:

$$Y^1 = \alpha_0^1 + \alpha_1^1 D + \alpha_2^1 X + \alpha_3^1 \theta + w_1$$

$$Y^2 = \alpha_0^2 + \alpha_1^2 D + \alpha_2^2 X + \alpha_3^2 \theta + w_2$$

$$Y^3 = \alpha_0^3 + \alpha_1^3 D + \alpha_2^3 X + \alpha_3^3 \theta + w_3$$

$$Y^4 = \alpha_0^4 + \alpha_1^4 D + \alpha_2^4 X + \alpha_3^4 \theta + w_4$$

where, $(\alpha_0^1, \alpha_0^2, \alpha_0^3, \alpha_0^4) = (2, 1, 1, 1)$, $(\alpha_2^1, \alpha_2^2, \alpha_2^3, \alpha_2^4) = (2, 1, 1, 1)$, $(\alpha_3^1, \alpha_3^2, \alpha_3^3, \alpha_3^4) = (2, 2, 2, 2)$. $(\alpha_1^1, \alpha_1^2, \alpha_1^3, \alpha_1^4)$, which represent the treatment effect of D on the outcomes Y^1, Y^2, Y^3 , and Y^4 , respectively, has been set equal to $(2, 3, 1, 1)$.

3.2 2SLS and SLF+IV

The specifications for D^* , X , D , z , θ , and Y are same as that of the baseline case. We again assume that $\begin{pmatrix} \theta \\ z \end{pmatrix} \sim \mathcal{N}\left(\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 1 & \sigma_{z\theta} \\ \sigma_{z\theta} & \sigma_z * \sigma_z \end{pmatrix}\right)$. Now, $\sigma_{z\theta} \in \{0, 0.025, 0.035, 0.045, 0.055, 0.065, 0.075\}$, where $\sigma_{z\theta} = 0$ is the uncontaminated IV case. Similarly, $\sigma_z \in \{0.16, 0.2, 0.24, 0.28, 0.32, 0.36\}$. Additionally, we have taken $\sigma_z = 1$ for the selected values of $\sigma_{z\theta}$, including the no contamination case, to make the simulations exercise in tandem with that of Banerjee and Basu (2021). So, the chosen values for $\sigma_{z\theta}$ and σ_z generate 49 combinations of strength and degree of contamination of the IV.

3.3 Estimation

We compare the OLS, 2SLS, SLF, and SLF+IV estimates based on percentage bias, standard error, coverage probability, and MSE, using Monte Carlo simulations across 1000 iterations of the dataset, each for $n = 500$ and $n = 1000$. The percentage bias has been calculated as:

$$\text{Percentage bias} = \frac{\text{Absolute Bias}}{\text{Actual Treatment Effect}} * 100$$

where the absolute bias has been calculated by taking the absolute difference between the actual values of $(\alpha^1, \alpha^2, \alpha^3, \alpha^4)$ and their estimated values obtained through simulations.

MSE has been calculated as:

$$\text{MSE} = \text{Bias}^2 + \text{Sampling variability}$$

Additionally, we have presented the first-stage F-statistic for the 2SLS estimate for a given combination of $\sigma_{z\theta}$ and σ_z . To decide whether to choose SLF or SLF+IV estimator based on the estimator bias, we have ascertained the thresholds for $\sigma_{z\theta}$ and σ_z from the simulations data. The simulations for SLF and SLF+IV employ the Maximum Likelihood Estimation (MLE) method to estimate the treatment effect. We have used the *gsem* package of Stata 18 (StataCorp, 2023) to obtain the simulations data.

3.4 Results and Discussion

Tables 1(a), 1(b), 1(c), and 1(d) pertain to Y^1, Y^2, Y^3 , and Y^4 , respectively, for $n = 500$. Results for $n = 1000$ have been depicted in tables A1(a), A1(b), A1(c), and A1(d) of the Appendix for the same four outcomes. In this section, the analyses thereafter pertain to all four outcome variables. The tables for $n = 500$ and $n = 1000$ show that the SLF estimator has a lower bias and MSE but a higher coverage probability than the OLS estimator. The SLF+IV estimate has either a lower coverage probability than the corresponding 2SLS estimate or both SLF+IV and 2SLS estimators have 0 coverage probability. Finally, the first-stage F-statistic of 2SLS estimates increases when $\sigma_{z\theta}$ increases for a given σ_z , as well as when σ_z increases for a given $\sigma_{z\theta}$.

The results for $n = 500$ and $n = 1000$ highlight that for a given σ_z , percentage bias and MSE of both 2SLS and SLF+IV estimates increase as $\sigma_{z\theta}$ increases. This shows that highly contaminated IVs are

Table 1(a): Y^1 , $n = 500$

True treatment effect = 2		Y ¹												0.075					
		Degrees of IV contamination \rightarrow 0 (no contamination)						0.025						0.035					
		Baseline cases			Naïve OLS			SLF			SLF+IV			2SLS			SLF+IV		
		std(z) = 0.16	Treatment effect estimate	3.763	2.608	1.185	2.003	5.145	2.263	5.802	2.736	6.430	4.379	7.017	5.991	7.321	6.326	7.589	6.447
			First stage F-stat			9.389		25.310		31.097		41.042		53.475		66.196		81.425	
			% Bias	88.164	30.416	40.770	0.160	157.232	13.174	190.097	36.776	221.507	118.931	250.839	199.552	266.055	216.314	279.468	222.333
			se	0.202	0.748	13.269	0.320	0.928	0.398	0.810	0.510	0.718	0.547	0.605	0.535	0.618	0.373	0.581	0.319
			Coverage Probability	0.000	0.701	0.963	0.958	0.093	0.949	0.003	0.858	0.000	0.443	0.000	0.064	0.000	0.002	0.000	0.000
			MSE	3.150	0.930	176.741	0.102	10.750	0.228	15.111	0.801	20.142	5.957	25.610	16.215	28.696	18.856	31.578	19.875
			std(z) = 0.2																
			Treatment effect estimate	3.763	2.608	1.758	2.005	4.358	2.227	4.958	2.396	5.420	2.675	5.896	3.660	6.468	5.207	6.786	5.991
			First stage F-stat			17.718		28.681		33.807		39.873		49.248		60.020		72.807	
			% Bias	88.164	30.416	12.100	0.263	117.905	11.368	147.887	19.807	171.005	33.736	82.996	223.425	160.369	239.284	199.565	
			se	0.202	0.748	1.587	0.318	0.865	0.382	0.769	0.449	0.700	0.541	0.646	0.613	0.608	0.555	0.573	0.456
			Coverage Probability	0.000	0.701	0.963	0.960	0.231	0.950	0.042	0.934	0.001	0.886	0.000	0.624	0.000	0.224	0.000	0.031
			MSE	3.150	0.930	2.578	0.101	6.310	0.198	9.339	0.358	12.188	0.748	15.599	3.131	20.337	10.595	23.231	16.138
			std(z) = 0.24																
			Treatment effect estimate	3.763	2.608	1.839	2.004	3.810	2.206	4.344	2.316	4.794	2.456	5.160	2.717	5.529	3.210	5.941	4.319
			First stage F-stat			23.496		38.947		46.094		52.295		61.557		60.614		70.722	
			% Bias	88.164	30.416	0.748	0.072	0.182	0.948	10.295	117.185	15.786	139.713	22.780	158.024	35.869	176.463	60.487	197.061
			se	0.202	0.748	1.251	0.313	0.813	0.366	0.730	0.413	0.673	0.480	0.630	0.567	0.594	0.705	0.563	0.708
			Coverage Probability	0.000	0.701	0.967	0.956	0.363	0.945	0.140	0.935	0.027	0.923	0.001	0.872	0.000	0.763	0.000	0.456
			MSE	3.150	0.930	1.592	0.098	3.936	0.176	6.026	0.270	8.261	0.438	10.385	0.836	12.808	1.960	15.850	5.881
			std(z) = 0.28																
			Treatment effect estimate	3.763	2.608	1.885	2.005	3.412	2.196	3.884	2.289	4.292	2.401	4.655	2.551	4.949	2.765	5.227	3.045
			First stage F-stat			28.041		36.563		42.879		49.331		56.344		63.956		72.211	
			% Bias	88.164	30.416	5.773	0.254	70.588	9.781	94.192	14.437	114.602	20.033	132.735	27.556	147.455	38.243	161.355	52.268
			se	0.202	0.748	1.046	0.308	0.760	0.353	0.694	0.384	0.644	0.433	0.606	0.573	0.603	0.549	0.731	0.779
			Coverage Probability	0.000	0.701	0.966	0.955	0.480	0.936	0.237	0.924	0.082	0.910	0.018	0.893	0.001	0.853	0.000	0.815
			MSE	3.150	0.930	1.112	0.095	2.570	0.163	4.631	0.231	5.669	0.348	7.415	0.564	9.026	0.949	10.712	1.626
			std(z) = 0.32																
			Treatment effect estimate	3.763	2.608	1.925	2.004	3.131	2.182	3.540	2.270	3.895	2.366	4.225	2.481	4.528	2.635	4.776	2.821
			First stage F-stat			35.367		44.475		48.305		57.288		63.798		69.136		71.027	
			% Bias	88.164	30.416	3.740	0.212	56.564	9.080	76.987	13.495	94.765	18.309	111.272	24.026	126.421	31.733	138.705	41.046
			se	0.202	0.748	0.907	0.303	0.711	0.340	0.659	0.363	0.617	0.396	0.584	0.443	0.556	0.520	0.532	0.601
			Coverage Probability	0.000	0.701	0.965	0.957	0.578	0.936	0.342	0.922	0.162	0.891	0.053	0.875	0.011	0.853	0.001	0.815
			MSE	3.150	0.930	0.828	0.092	1.785	0.149	2.805	0.205	3.973	0.291	5.294	0.427	7.672	0.674	7.989	0.105
			std(z) = 0.36																
			Treatment effect estimate	3.763	2.608	1.941	2.006	2.934	2.176	3.278	2.254	3.593	2.345	3.886	2.442	4.154	2.554	4.407	2.697
			First stage F-stat			55.211		70.412		73.357		77.091		80.454		85.273		93.186	
			% Bias	88.164	30.416	2.967	0.305	46.677	8.785	63.889	12.695	79.635	17.322	94.298	22.093	107.698	27.684	120.339	34.855
			se	0.202	0.748	0.807	0.298	0.664	0.630	0.932	0.429	0.909	0.245	0.888	0.107	0.860	0.036	0.838	0.006
			Coverage Probability	0.000	0.701	0.965	0.953	0.655	0.809	1.312	0.139	0.202	0.186	2.884	0.258	3.872	0.360	4.929	0.512
			MSE	3.150	0.930	0.133	0.057	546.090	0.293	8.800	4.210	12.262	5.698	15.899	2.180	2.456	2.211	2.528	2.247
			std(z) = 1																
			Treatment effect estimate	3.763	2.608	1.992	2.006	2.176	2.084	2.245	2.114	2.318	2.146	2.389	2.180	2.456	2.211	2.528	2.247
			First stage F-stat			88.164		30.416		40.403		42.939		45.883		49.144		53.534	
			% Bias	88.164	30.416	0.202	0.365	0.239	0.359	0.243	0.357	0.245	0.355	0.249	0.333	0.249	0.351	0.251	0.348
			se	0.000	0.701	0.956	0.942	0.922	0.924	0.888	0.915	0.846	0.897	0.775	0.877	0.712	0.858	0.631	0.827
			Coverage Probability	0.000	0.701	0.956	0.942	0.930	0.133	0.057	0.160	0.066	0.187	0.073	0.227	0.094	0.276	0.082	0.125

Notes:

First stage F-stat: F statistic for the first stage regression of 2SLS

% bias: percentage bias

se: standard error

MSE: Mean Squared Error

Table 1(b): $Y^2, n = 500$

True treatment effect = 3		Degrees of IV contamination \rightarrow 0 (no contamination)										0.025		0.035		0.045		0.055		0.065		0.075	
IV Strength \downarrow	std(z) = 0.16	Baseline cases		Naive OLS		SLF		2SLS		SLF+IV		2SLS		SLF+IV		2SLS		SLF+IV		2SLS		SLF+IV	
std(z) = 0.2	Treatment effect estimate	4.752	3.593	2.211	3.007	6.123	3.269	6.807	3.742	7.418	5.378	8.001	6.985	8.295	7.319	8.564	7.440						
	First stage F-stat	-	9.389	-	25.310	-	31.097	-	41.042	-	53.475	-	66.196	-	81.425	-							
	% Bias	58.387	19.757	26.316	0.237	104.088	8.958	126.888	24.730	147.279	79.271	166.703	132.843	143.972	185.452	147.990	-						
	se	0.202	0.747	11.794	0.320	0.926	0.397	0.811	0.510	0.718	0.545	0.663	0.532	0.616	0.372	0.579	0.319						
	Coverage Probability	0.000	0.705	0.964	0.969	0.089	0.952	0.002	0.862	0.000	0.445	0.000	0.062	0.000	0.001	0.000	0.000	0.000					
	MSE	3.109	0.909	139.722	0.102	10.609	0.230	15.148	0.811	20.037	5.933	25.451	16.166	28.418	18.794	31.289	19.813						
std(z) = 0.24	Treatment effect estimate	4.752	3.593	2.734	3.010	5.333	3.232	5.948	3.400	6.414	3.677	6.873	4.651	7.456	6.204	7.761	6.984						
	First stage F-stat	-	17.718	-	8.852	-	3.42	28.681	-	33.807	-	13.342	113.733	22.561	129.097	55.030	148.541	106.812	72.807	-			
	% Bias	58.387	19.757	0.202	0.747	1.584	0.317	0.866	0.381	0.769	0.448	0.700	0.541	0.647	0.612	0.609	0.573	0.455					
	se	0.000	0.705	0.966	0.962	0.233	0.946	0.036	0.927	0.001	0.893	0.000	0.616	0.000	0.222	0.000	0.000	0.031					
	Coverage Probability	0.000	0.909	2.581	0.101	6.193	0.199	9.282	0.361	12.144	0.751	15.418	3.100	20.229	10.574	22.995	16.078						
	MSE	3.109	0.909	2.814	3.005	4.790	3.208	5.324	3.319	5.771	3.458	6.141	3.722	6.495	4.211	6.910	5.305						
std(z) = 0.28	Treatment effect estimate	4.752	3.593	23.496	-	38.947	-	46.094	-	46.094	-	52.295	-	61.557	-	60.614	-	70.722	-				
	First stage F-stat	-	19.757	6.211	0.156	59.671	0.515	59.671	77.477	10.621	92.355	15.262	104.707	24.074	116.488	40.379	130.331	76.822					
	% Bias	58.387	19.757	0.202	0.747	1.251	0.313	0.813	0.365	0.730	0.412	0.673	0.478	0.629	0.565	0.595	0.703	0.563					
	se	0.000	0.705	0.967	0.954	0.367	0.943	0.136	0.935	0.024	0.921	0.001	0.867	0.000	0.753	0.000	0.000	0.446					
	Coverage Probability	0.000	0.909	1.601	0.098	3.866	0.177	5.936	0.271	8.129	0.438	10.263	0.841	12.566	1.962	15.604	5.812						
	MSE	3.109	0.909	2.862	3.006	4.390	3.199	4.860	3.290	5.261	3.401	5.628	3.553	5.929	3.767	6.202	4.052						
std(z) = 0.32	Treatment effect estimate	4.752	3.593	28.041	-	36.563	-	42.879	-	49.331	-	56.344	-	63.056	-	63.056	-	72.211	-				
	First stage F-stat	-	19.757	4.605	0.185	46.337	0.635	62.001	9.683	75.376	13.371	87.610	18.427	97.641	25.564	106.720	35.067						
	% Bias	58.387	19.757	0.202	0.747	0.104	0.308	0.760	0.353	0.694	0.383	0.644	0.432	0.605	0.508	0.573	0.601	0.548					
	se	0.000	0.705	0.968	0.950	0.498	0.930	0.241	0.919	0.079	0.916	0.016	0.888	0.001	0.854	0.000	0.000	0.788					
	Coverage Probability	0.000	0.909	1.119	0.095	2.510	0.164	3.941	0.231	5.528	0.348	7.074	0.563	8.908	0.949	10.551	1.636						
	MSE	3.109	0.909	2.899	3.006	4.110	3.184	4.519	3.205	4.874	3.274	5.202	3.368	5.202	3.482	5.206	3.635	5.755					
std(z) = 0.36	Treatment effect estimate	4.752	3.593	35.367	3.383	0.198	36.989	0.136	50.646	9.128	62.471	12.272	73.414	16.054	83.518	21.177	91.823	27.338					
	First stage F-stat	-	19.757	0.747	0.908	0.302	0.711	0.339	0.658	0.362	0.617	0.395	0.583	0.441	0.555	0.519	0.532	0.599					
	% Bias	58.387	19.757	0.202	0.747	0.964	0.950	0.956	0.929	0.342	0.909	0.166	0.896	0.055	0.883	0.013	0.858	0.001	0.804				
	se	0.000	0.705	0.909	0.834	0.091	1.737	0.149	2.742	0.206	3.893	0.292	5.191	0.427	5.86	0.673	7.871	1.032					
	Coverage Probability	0.000	0.909	0.659	0.089	1.272	0.139	1.975	0.185	0.824	0.257	3.798	0.358	4.839	0.511	5.958	0.511						
	MSE	3.109	0.909	3.593	2.980	2.996	3.164	3.074	3.235	3.105	3.308	3.139	3.377	3.171	3.445	3.202	3.517	3.237					
std(z) = 1	Treatment effect estimate	4.752	3.593	54.090	-	57.012	-	57.265	-	569.392	-	589.101	-	598.454	-	613.645	-						
	First stage F-stat	-	19.757	0.653	0.141	5.459	2.475	7.831	3.516	10.269	4.622	12.566	5.688	14.847	6.741	17.247	-						
	% Bias	58.387	19.757	0.202	0.747	0.328	0.359	0.243	0.356	0.245	0.347	0.354	0.352	0.247	0.350	0.251	0.348	0.253					
	se	0.000	0.705	0.959	0.938	0.934	0.937	0.902	0.929	0.851	0.912	0.799	0.891	0.727	0.870	0.664	0.841	0.759					
	Coverage Probability	0.000	0.909	0.133	0.057	0.156	0.065	0.182	0.071	0.221	0.080	0.266	0.091	0.321	0.104	0.389	0.120						
	MSE	3.109	0.909	2.980	2.996	3.164	3.074	3.235	3.105	3.308	3.139	3.377	3.171	3.445	3.202	3.517	3.237						

Notes:

First stage F-stat: F statistic for the first stage regression of 2SLS

% bias: percentage bias

se: standard error

MSE: Mean Squared Error

std(z): standard deviation of z

Table 1(c): $Y^3, n = 500$

Notes:

First stage F-stat: F statistic for the first stage regression of 2SLS

% bias: percentage bias

use: standard error

MSE: Mean Squared Error

Table 1(d): $Y^4, n = 500$

Notes:

First stage F-stat: F statistic for the first stage regression of 2SLS

% bias: percentage bias

use: standard error of treatment effect

MSE: Mean Squared Error

$\text{std}(z)$: standard deviation of z

associated with higher 2SLS and SLF+IV estimator bias for any given IV strength. On the other hand, for a given non-zero $\sigma_{z\theta}$, percentage bias and MSE of both 2SLS and SLF+IV estimates decline as the σ_z increases for both sample sizes. This shows that stronger IVs are associated with lower 2SLS and SLF+IV estimator bias. However, this virtue does not hold when the IV is uncontaminated. Hence, by combining these conjectures, we can say that 2SLS and SLF+IV estimates have lower (greater) biases when the IV is strong (weak) and less (more) contaminated.

As described in our theoretical framework, there exists a threshold $\sigma_{z\theta}^*$ for a given strength of the IV, such that when $\sigma_{z\theta}$ is lower (higher) than that threshold, the SLF+IV estimate will have a lower (greater) bias than the SLF estimate. Our simulation results are consistent with the theoretical framework. For instance, when $n = 500$ and $\sigma_z = 0.16$, $\sigma_{z\theta}^*$ lies between 0.025 and 0.035, because when $\sigma_z = 0.16$, the bias of the SLF+IV estimator is greater (lower) than that of the SLF estimator when $\sigma_{z\theta} \leq 0.025$ ($\sigma_{z\theta} \geq 0.035$). In the same manner, when $n = 500$, $\sigma_{z\theta}^*$ lies between 0.035 and 0.045, 0.045 and 0.055, 0.055 and 0.065, 0.055 and 0.065, and 0.075, when the given σ_z is 0.2, 0.24, 0.28, 0.32, and 0.36, respectively. Similarly, when $n = 1000$, $\sigma_{z\theta}^*$ lies between 0.025 and 0.035, 0.035 and 0.045, 0.035 and 0.045, 0.045 and 0.055, 0.045 and 0.055, and 0.045 and 0.055, when the given σ_z is 0.16, 0.2, 0.24, 0.28, 0.32, and 0.36, respectively. Note that when $\sigma_z = 1$, $\sigma_{z\theta}^*$ lies above 0.075, for both $n = 500$ and $n = 1000$. Hence, the threshold $\sigma_{z\theta}^*$ determines the choice between SLF and SLF+IV estimator for a given σ_z and sample size.

Our theoretical framework also elucidates that there exists a threshold σ_z^* for a given $\sigma_{z\theta}$, such that when σ_z is lower (higher) than σ_z^* , the SLF+IV estimate will have a greater (lower) bias than the SLF estimate. Our simulation results are again consistent with this theoretical exposition. For instance, when $n = 500$ and $\sigma_{z\theta} = 0.045$, σ_z^* lies between 0.2 and 0.24, because when $\sigma_{z\theta} = 0.045$, the bias of the SLF+IV estimator is lower (greater) than that of the SLF estimator when $\sigma_z \leq 0.2$ ($\sigma_z \geq 0.24$). In the same manner, when $n = 500$, σ_z^* lies below 0.025, between 0.16 and 0.2, 0.24 and 0.28, 0.32 and 0.36, and 0.36 and 1, when $\sigma_{z\theta}$ is 0.025, 0.035, 0.055, 0.065, and 0.075, respectively. When $n = 1000$, σ_z^* lies below 0.025, between 0.16 and 0.2, between 0.24 and 0.28, between 0.24 and 0.28, between 0.36 and 1, between 0.36 and 1, and between 0.36 and 1, when $\sigma_{z\theta}$ is 0.025, 0.035, 0.045, 0.055, 0.065, and 0.075, respectively. Hence, the threshold σ_z^* determines the choice between SLF and SLF+IV estimator for a given $\sigma_{z\theta}$ and sample size.

So, by combining the attributes of thresholds $\sigma_{z\theta}^*$ and σ_z^* simultaneously, we find that the SLF+IV estimator is preferred over the SLF estimator when the IV is strong and less contaminated, while the SLF estimator is preferred over the SLF+IV estimator when the IV is weak and highly contaminated. Our simulation results unveil that the strength and degree of contamination of the IV have opposite effects on the bias of the SLF+IV estimator. So, when the IV strength and degree of contamination increase (decrease) simultaneously, then the SLF+IV estimator remains a better (worse) choice than the SLF estimator if the effect of strengthening (weakening) of the IV *outweighs* the effect of rise (fall) in the degree of contamination of the IV.

Figures A1-A4 in the Appendix are the heatplots for outcomes Y^1, Y^2, Y^3, Y^4 , respectively, for $n = 500$, wherein each block represents the bias of the SLF+IV estimator for a given combination of the strength of the IV and degree of contamination. The x-axis represents the strength of the IV, and the y-axis represents the degree of contamination. Note that we have considered only $\sigma_{z\theta} \in \{0.025, 0.035, 0.045, 0.055, 0.065, 0.075\}$, and $\sigma_z \in \{0.16, 0.2, 0.24, 0.28, 0.32, 0.36\}$ for the construction of the heatplots. The intensity of the orange colour represents the magnitude of the percentage bias of the SLF+IV estimators for Y^1, Y^2, Y^3, Y^4 , respectively, such that darker shades of orange represent higher biases of the estimator. In the legends of each figure, we have mentioned the extreme values of the percentage bias of the SLF+IV estimator for each outcome on the right side of the colour intensity bar, while the corresponding bias of the SLF estimator is mentioned on the left side of the bar. Note that blocks with lower IV strength and higher degree of contamination have darker shades

of orange. We have displayed the preferred estimator between SLF and SLF+IV on each block after comparing the bias of the SLF+IV estimator with that of the SLF estimator. The heatplots show that the SLF+IV (SLF) estimator is preferred over the SLF (SLF+IV) estimator at blocks with higher (lower) IV strength and lower (higher) degree of contamination.

The simulation results reveal regarding $\sigma_{z\theta}^*$ and σ_z^* that $\sigma_{z\theta}^*(n = 500, \bar{\sigma}_z) \geq \sigma_{z\theta}^*(n = 1000, \bar{\sigma}_z)$, and $\sigma_z^*(n = 500, \bar{\sigma}_{z\theta}) \leq \sigma_z^*(n = 1000, \bar{\sigma}_{z\theta})$. These results indicate that as the sample size increases, the threshold strength and degree of contamination of the IV for choosing between SLF and SLF+IV estimator declines and increases, respectively. This indicates that larger sample sizes are associated with SLF being a better choice than SLF+IV for a given strength and degree of contamination of the IV. So, the sample size creates a trade-off between SLF and SLF+IV estimators, wherein SLF+IV is a better choice for smaller sample sizes, while SLF is a better choice for larger sample sizes.

4 Empirical Application: Impact of Obesity on Different Indicators of the Health Status

4.1 Introduction

Obesity is a prevalent, complex, progressive and relapsing chronic disease characterized by abnormal or excessive body fat that impairs health (Wharton *et al.*, 2020). It has become a serious public health concern in recent years. A person is said to be obese if her Body Mass Index (BMI) is above 30 kg/m^2 (World Health Organization, 1998). Obesity is associated with many non-communicable diseases—including hypertension and diabetes mellitus (Segula, 2014). Hypertension is a frequent, chronic, age-related disorder which often entails debilitating cardiovascular and renal complications. The diastolic and systolic readings of blood pressure are key indicators of hypertension (Staessen, 2003). Diabetes mellitus is a chronic disorder of glucose metabolism, which leads to microvascular (retinopathy, nephropathy, neuropathy) and macrovascular (ischaemic heart disease, stroke, peripheral vascular disease) complications (Forouhi and Wareham, 2010). Diabetes mellitus is associated with a high blood sugar level. Apart from hypertension and diabetes, studies have shown that there is a strong association between obesity and anaemia (Saad and Qutob, 2022; Cepeda-Lopez and Baye, 2020; Jordaan *et al.*, e2020). Experts have opined that obesity creates disturbances at the endothelial, hormonal, and inflammatory levels, which potentially leads to anaemic state (Saad and Qutob, 2022). An individual is diagnosed to be anaemic or non-anaemic based on the haemoglobin count in her blood.

Our empirical application aims to test the causal impact of obesity on the risk of hypertension, diabetes mellitus and anaemia among women in India. It is, however, difficult to obtain an unbiased estimate of the impact due to the presence of unobservable confounders that affect both the prevalence of obesity in a woman and the risk of hypertension, diabetes, and anaemia. For instance, genetic factors can potentially affect the prevalence of obesity and indicators for those diseases. In that case, an IV may be employed to create some exogenous variation in the treatment variable to mitigate the impact of any unobservable confounders. We have considered the *arm circumference* as the IV. The IV is likely to be a strong IV, as several studies have revealed that arm circumference is strongly associated with obesity (Kiran *et al.*, 2022; Dereje *et al.*, 2022; Craig *et al.*, 2014). However, our IV is potentially contaminated. The arm circumference can affect the cuff size for blood pressure measurement, which can affect the blood pressure reading of an individual (Loenekke *et al.*, 2016). A larger arm circumference is potentially linked with reduced insulin resistance, which raises the risk of a high blood glucose level. Such insulin resistance can be independent of the BMI level of an individual (Wang *et al.*, 2023). Moreover, a larger arm circumference is strongly associated with the haemoglobin level of an individual (Ahankari *et al.*, 2020), although it is not clear whether the causal impact of arm circumference on the haemoglobin level is through the obesity of an individual or due

to any other unaccounted factor. Hence, in view of the possibility of the IV being contaminated, a highly contaminated IV shall make the 2SLS estimator more biased than the OLS estimator, even if the OLS estimator itself is highly biased.

We motivate the application of SLF and SLF+IV methods, along with the traditional estimators (OLS and 2SLS). Both the latent factor estimators have the virtue of being less biased than the traditional ones, which has been elucidated in the preceding sections. Using the empirical application, we also aim to demonstrate the virtue of the latent factor models in capturing the effects of omitted variable(s) that are correlated with treatment and outcome measurements.

4.2 Data

The data for this analysis comes from the latest fifth Round of the National Family Health Survey (NFHS-5), conducted in 2019-21 (International Institute for Population Sciences (IIPS) & ICF, 2021). The Indian NFHS is comparable to the Demographic Health Surveys (DHS) conducted in many other countries. NFHS is a nationally representative repeated cross-sectional study, which was started in 1992-93 (NFHS-1), and the latest round was conducted in 2019-21 (NFHS-5). The NFHS series furnishes data on socioeconomic, demographic, health, and nutritional indicators across 707 districts of India, covering all the states and union territories. This data is delineated by demographic characteristics at both national and state levels.

Our dataset contains women aged between 30 and 40 years. We have considered only those respondents whose BMI lies in either the healthy weight category ($18.5 < \text{BMI} < 24.9$) or the obese category ($\text{BMI} > 30$), wherein obese women constitute the treatment group, while healthy weight women constitute the control group. So, our treatment variable is a binary variable, which takes the value 1 if the woman is obese, and 0 if the woman has a healthy weight. We have taken arm circumference as a continuous IV. The final sample size, thus, stands at 65,034.

We have considered four continuous outcome variables: (i) log of average diastolic reading, (ii) log of average systolic reading, (iii) log of blood glucose level, and (iv) log of haemoglobin level. The outcome variable (i) has been generated by taking the logarithm of the average of the three diastolic readings (mmHg) of the blood pressure of the woman. Similarly, the outcome variable (ii) has been generated by taking the logarithm of the average of the three systolic readings (mmHg) of the blood pressure of the woman. The outcome variable (iii) has been generated by taking the logarithm of the blood sugar level (mg/dL) of the woman. Finally, the outcome variable (iv) has been generated by taking the logarithm of the haemoglobin level (g/dL) of the woman.

The control variables include the place of residence, current age, alcohol consumption and the dietary habits of the woman. The place of residence has been categorized into rural (reference category) and urban. The current age of the woman is a continuous variable. The alcohol consumption of the woman has been categorized into whether the woman consumes alcohol or not (reference category). We have controlled the consumption of milk/curd, pulses/beans, dark green leafy vegetables, fruits, fish, chicken/meat, fried food, and aerated drinks. We have categorized the consumption of any of the food items into no consumption (reference category), daily consumption, weekly consumption, or occasional consumption.

4.3 Methods

We have studied the difference in the means of the outcome variables, IV, and control variables between the treatment group (obese women) and the control group (healthy-weight women) using a t-test (in case of continuous variable) and Chi-square test (in case of binary variable).

In the baseline cases, we have considered the traditional models: OLS and 2SLS. In the OLS model, we have regressed each outcome variable on the treatment indicator and all the control variables. We have used the *regress* command of Stata 18 to run the OLS models. The 2SLS model comprises two regressions in two stages. In the first stage regression, we regressed the treatment variable on the IV and all the other control variables and then regressed each outcome variable on the treatment variable estimated from the first stage regression and all the other control variables. We have used the *ivregress* command of Stata 18 to run the 2SLS model. The standard errors obtained for both OLS and 2SLS estimates are heteroscedasticity robust.

Next, we estimated SLF and SLF+IV models with the same four outcome, treatment, and control variables. The models have been estimated using the *gsem* package of Stata 18 (StataCorp, 2023). The structural parameters of both SLF and SLF+IV models will be identified as the number of outcomes chosen by us fulfils the identification criteria for both SLF and SLF+IV models. All four models are estimated using the technique of Maximum Likelihood Estimation (MLE). Note that we have assumed that the treatment variable for both SLF and SLF+IV have been drawn from a probit model, in order to make our model similar to our simulations exercise models. The standard errors obtained for both SLF and SLF+IV estimates are heteroscedasticity robust.

We have tabulated the summary statistics that outline the means and standard deviations of all variables between the treatment and control groups and the p-values of the differences in means. We have tabulated the treatment effect estimates of obesity on the log of average diastolic reading, log of average systolic reading, log of blood glucose level, and log of haemoglobin level, along with the standard deviation, number of observations and mean of observation. We have also shown the treatment effects graphically, along with 95 % confidence intervals.

4.4 Results

4.4.1 Summary Statistics

In Table 2, the summary statistics reveal that all four outcome variables and the IV have significantly greater means in the treatment group than in the control group. The mean age of women in the treatment group is significantly greater than that in the control group. The women in both the treatment and control groups majorly belong to rural areas, wherein the treatment group has a larger proportion of urban women than the control group. The women in both groups are mostly non-consumers of alcohol, wherein the treatment group has a lower proportion of alcohol consumers than the control group. Most of the respondents in the treatment group are daily consumers of milk/curd, daily consumers of pulses/beans, daily consumers of dark green leafy vegetables, weekly consumers of fruits, non-consumers of fish, weekly consumers of chicken/meat, occasional consumers of fried food, and occasional consumers of aerated drinks. On the other hand, most of the respondents in the control group are daily consumers of milk/curd, daily consumers of pulses/beans, daily consumers of dark green leafy vegetables, occasional consumers of fruits, weekly consumers of fish, occasional consumers of chicken/meat, occasional consumers of fried food, and occasional consumers of aerated drinks. The treatment group has significantly greater proportions of daily consumers of milk/curd, daily consumers of pulses/beans, daily and weekly consumers of fruits, non-consumers of fish, daily consumers of fish, non-consumers and daily consumers of chicken/meat, non-consumers and weekly of fried food, and weekly consumers of aerated drinks than the control group. On the other hand, the control group has significantly greater proportions of weekly and occasional milk/curd consumers, non-consumers and occasional consumers of pulses/beans, non-consumers and occasional consumers of fruits, weekly and occasional consumers of fish, occasional consumers of chicken/meat, daily and occasional consumers of fried food, and non-consumers and occasional consumers of aerated drinks than the treatment group.

Table 2: Summary Statistics

VARIABLES	Treatment group: Obese women Mean [a]	SD	Control group: Healthy-weight women Mean [b]	SD	p-value of difference [†] [a] - [b]
Outcomes					
Log of average diastolic reading	4.404 (N = 3,511)	0.109	4.349 (N = 35,594)	0.110	< 0.001
Log of average systolic reading	4.775 (N = 3,511)	0.113	4.751 (N = 35,595)	0.108	< 0.001
Log of blood glucose level	4.739 (N = 3,640)	0.239	4.679 (N = 37,902)	0.171	< 0.001
Log of haemoglobin count	4.763 (N = 3,640)	0.135	4.737 (N = 37,867)	0.154	< 0.001
Instrumental variable					
Arm circumference	29.552 (N = 3,728)	3.349	24.342 (N = 38,279)	2.310	< 0.001
Control variables					
Place of residence:					
Urban	0.437 (N = 3,740)	0.008	0.214 (N = 38,380)	0.002	< 0.001
Rural [‡]	0.563 (N = 3,740)	0.008	0.786 (N = 38,380)	0.002	< 0.001
Current age	33.294 (N = 3,740)	2.805	32.988 (N = 38,380)	2.824	< 0.001
Alcohol consumption:					
Yes	0.014 (N = 3,740)	0.002	0.024 (N = 38,380)	0.001	< 0.001
No [§]	0.986 (N = 3,740)	0.002	0.976 (N = 38,380)	0.001	< 0.001
Milk/Curd consumption:					
Never [¶]	0.056 (N = 3,740)	0.004	0.076 (N = 38,380)	0.001	< 0.001
Daily	0.568 (N = 3,740)	0.008	0.407 (N = 38,380)	0.003	< 0.001
Weekly	0.206 (N = 3,740)	0.007	0.253 (N = 38,380)	0.002	< 0.001
Occasionally	0.170 (N = 3,740)	0.006	0.263 (N = 38,380)	0.002	< 0.001
Pulses/Beans consumption:					
Never [¶]	0.003 (N = 3,740)	0.001	0.003 (N = 38,380)	0.000	< 0.001
Daily	0.497 (N = 3,740)	0.008	0.456 (N = 38,380)	0.003	< 0.001
Weekly	0.428 (N = 3,740)	0.008	0.438 (N = 38,380)	0.003	0.222
Occasionally	0.072 (N = 3,740)	0.004	0.103 (N = 38,380)	0.002	< 0.001
Dark green leafy vegetables consumption:					
Never [¶]	0.001 (N = 3,740)	0.000	0.002 (N = 38,380)	0.000	0.105
Daily	0.546 (N = 3,740)	0.008	0.549 (N = 38,380)	0.003	0.712
Weekly	0.367 (N = 3,740)	0.008	0.358 (N = 38,380)	0.002	0.313
Occasionally	0.087 (N = 3,740)	0.005	0.091 (N = 38,380)	0.001	0.409
Fruits consumption:					
Never [¶]	0.011 (N = 3,740)	0.002	0.014 (N = 38,380)	0.001	0.086
Daily	0.240 (N = 3,740)	0.007	0.112 (N = 38,380)	0.002	< 0.001
Weekly	0.414 (N = 3,740)	0.008	0.365 (N = 38,380)	0.002	< 0.001
Occasionally	0.338 (N = 3,740)	0.008	0.515 (N = 38,380)	0.003	< 0.001
Fish consumption:					
Never [¶]	0.378 (N = 3,740)	0.008	0.280 (N = 38,380)	0.002	< 0.001
Daily	0.068 (N = 3,740)	0.004	0.042 (N = 38,380)	0.001	< 0.001
Weekly	0.285 (N = 3,740)	0.007	0.205 (N = 38,380)	0.002	0.010
Occasionally	0.270 (N = 3,740)	0.007	0.374 (N = 38,380)	0.002	< 0.001
Chicken/meat consumption:					
Never [¶]	0.340 (N = 3,740)	0.008	0.254 (N = 38,380)	0.002	< 0.001
Daily	0.060 (N = 3,740)	0.003	0.017 (N = 38,380)	0.001	< 0.001
Weekly	0.344 (N = 3,740)	0.008	0.311 (N = 38,380)	0.002	0.691
Occasionally	0.291 (N = 3,740)	0.007	0.389 (N = 38,380)	0.002	< 0.001
Fried food consumption:					
Never [¶]	0.051 (N = 3,740)	0.004	0.042 (N = 38,380)	0.001	0.017
Daily	0.089 (N = 3,740)	0.005	0.102 (N = 38,380)	0.002	0.014
Weekly	0.370 (N = 3,740)	0.008	0.348 (N = 38,380)	0.002	0.007
Occasionally	0.490 (N = 3,740)	0.008	0.507 (N = 38,380)	0.003	0.040
Aerated drinks consumption:					
Never [¶]	0.155 (N = 3,740)	0.006	0.175 (N = 38,380)	0.002	0.003
Daily	0.030 (N = 3,740)	0.003	0.030 (N = 38,380)	0.001	0.925
Weekly	0.152 (N = 3,740)	0.006	0.130 (N = 38,380)	0.002	< 0.001
Occasionally	0.662 (N = 3,740)	0.008	0.665 (N = 38,380)	0.002	0.740

Notes:

N: number of observations after adjusting for missing values

SD: standard deviation

†p-values reported for Chi-square test (in case of binary variables) and t-test (in case of continuous variables)

‡reference category

4.4.2 Strength of the IV

The IV has a significant effect on the treatment variable ($p < 0.001$) after controlling for the control variables. The F-statistic of the first stage regression of the 2SLS model is 255.79, which is much larger than 10. So, as per the thumb rule proposed by Staiger and Stock (1997), our IV is a strong IV.

4.4.3 Effect of obesity on the log of average diastolic reading

In Table 3, all four estimators reveal that obesity leads to a significant rise in average diastolic reading. The coefficients of the shared latent factor are significant in both SLF and SLF+IV models, which shows that there are potential omitted variables in both OLS and 2SLS models. Table A2 in the Appendix shows that staying in urban areas and daily consumption of fish are significantly associated with a lower diastolic reading, while alcohol consumption, daily consumption of pulses/beans, consumption of chicken/meat, and daily consumption of aerated drinks are significantly associated with higher diastolic reading. Older women are more likely to have a higher diastolic reading. In Figure 1, the OLS and SLF treatment estimates are close to each other and have relatively narrower 95% confidence intervals than 2SLS and SLF+IV estimates. The 2SLS estimate has the largest magnitude and the widest confidence interval.

4.4.4 Effect of obesity on the log of average systolic reading

In Table 4, all four estimators reveal that obesity leads to a significant rise in average systolic reading. The coefficients of the shared latent factor are significant in both SLF and SLF+IV models, which shows that there are potential omitted variables in both OLS and 2SLS models. Table A2 in the Appendix shows that staying in urban areas, daily consumption of fish, and daily consumption of fruits are significantly associated with a lower diastolic reading, while alcohol consumption, consumption of milk/curd, occasional consumption of pulses/bean, and consumption of chicken/meat are significantly associated with higher systolic reading. Older women are more likely to have a higher systolic reading. In Figure 2, the OLS and SLF treatment estimates are close to each other and have relatively narrower 95% confidence intervals than 2SLS and SLF+IV estimates. The 2SLS estimate has the largest magnitude and the widest confidence interval.

4.4.5 Effect of obesity on the log of blood glucose level

In Table 5, all four estimators reveal that obesity leads to a significant rise in blood glucose level. The coefficients of the shared latent factor are significant in both SLF and SLF+IV models, which shows that there are potential omitted variables in both OLS and 2SLS models. Table A2 in the Appendix shows that chicken/meat consumption is significantly associated with a lower blood glucose level, while daily and weekly consumption of milk/curd, consumption of fish, and consumption of aerated drinks are significantly associated with a higher blood glucose level. Older women are more likely to have a higher systolic reading. In Figure 3, the OLS, SLF, and SLF+IV are almost close to each other and almost have the same 95% confidence interval width, while the 2SLS estimate has the largest magnitude and the widest confidence interval.

4.4.6 Effect of obesity on the log of haemoglobin level

In Table 6, all four estimators reveal that obesity leads to a significant rise in haemoglobin level in the blood. The coefficients of the shared latent factor are significant in both SLF and SLF+IV models, which shows that there are potential omitted variables in both OLS and 2SLS models. Table A2 in the Appendix shows that occasional consumption of milk/curd is significantly associated with a lower haemoglobin level in the blood while staying in urban areas, daily consumption of milk/curd, occasional consumption of pulses/beans, daily and weekly consumption of fruits, daily and occasional

Table 3: Effect of obesity on log of average diastolic reading

	log of average diastolic reading			
	OLS	2SLS	SLF	SLF+IV
Treatment effect: obesity [†]	0.0570*** (0.00194)	0.105*** (0.00398)	0.0577*** (0.00194)	0.0841*** (0.00320)
Number of Observations	39,105	39,101	42,120	42,074
Mean of the outcome	4.355	4.355	4.355	4.355
Coefficient of the latent factor			-21.00*** (0.880)	4.214*** (0.0244)

Notes:

† after controlling for age, place of residence, alcohol consumption and dietary habits
*** $p < 0.01$, ** $p < 0.05$, * $p < 0.1$

Robust standard errors in the parentheses

Table 4: Effect of obesity on log of average systolic reading

	log of average systolic reading			
	OLS	2SLS	SLF	SLF+IV
Treatment effect: obesity	0.0270*** (0.00199)	0.0942*** (0.00415)	0.0276*** (0.00199)	0.0505*** (0.00308)
Robust standard error				
Number of Observations	39,106	39,102	42,120	42,074
Mean of the outcome	4.751	4.751	4.751	4.751
Coefficient of the latent factor			-17.86*** (0.815)	4.626*** (0.0204)

Notes:

† after controlling for age, place of residence, alcohol consumption and dietary habits
** $p < 0.01$, ** $p < 0.05$, * $p < 0.1$

Robust standard errors in the parentheses

Table 5: Effect of obesity on log of blood glucose level

	log of blood glucose level			
	OLS	2SLS	SLF	SLF+IV
Treatment effect: obesity [†]	0.0583*** (0.00407)	0.0978*** (0.00713)	0.0585*** (0.00408)	0.0604*** (0.00413)
Robust standard error				
Number of Observations	41,542	41,475	42,120	42,074
Mean of the outcome	4.686	4.686	4.686	4.686
Coefficient of the latent factor			-1.437*** (0.309)	4.571*** (0.0240)

Notes:

† after controlling for age, place of residence, alcohol consumption and dietary habits
** $p < 0.01$, ** $p < 0.05$, * $p < 0.1$

Robust standard error in the parentheses

Table 6: Effect of obesity on haemoglobin level

	log of blood glucose level			
	OLS	2SLS	SLF	SLF+IV
Treatment effect: obesity [†]	0.0205*** (0.00243)	0.121*** (0.00554)	0.0213*** (0.00243)	0.0279*** (0.00254)
Robust standard error				
Number of Observations	41,507	41,441	42,120	42,074
Mean of the outcome	4.738	4.738	4.738	4.738
Coefficient of the latent factor			-5.145*** (0.280)	4.658*** (0.0299)

Notes:

† after controlling for age, place of residence, alcohol consumption and dietary habits

** * $p < 0.01$, ** $p < 0.05$, * $p < 0.1$

Robust standard error in the parentheses

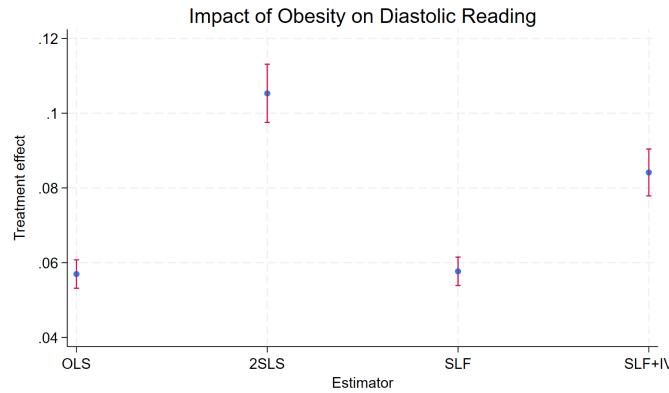


Fig 1: OLS, 2SLS, SLF, and SLF+IV estimates of the effect of obesity on diastolic reading, with 95% confidence intervals

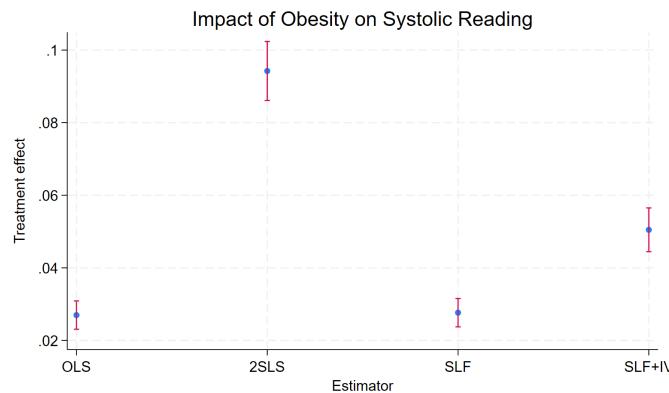


Fig 2: OLS, 2SLS, SLF, and SLF+IV estimates of the effect of obesity on systolic reading, with 95% confidence intervals

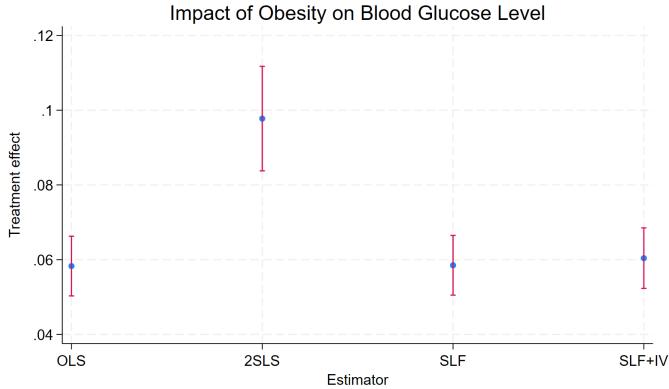


Fig 3: OLS, 2SLS, SLF, and SLF+IV estimates of the effect of obesity on blood glucose level, with 95% confidence intervals

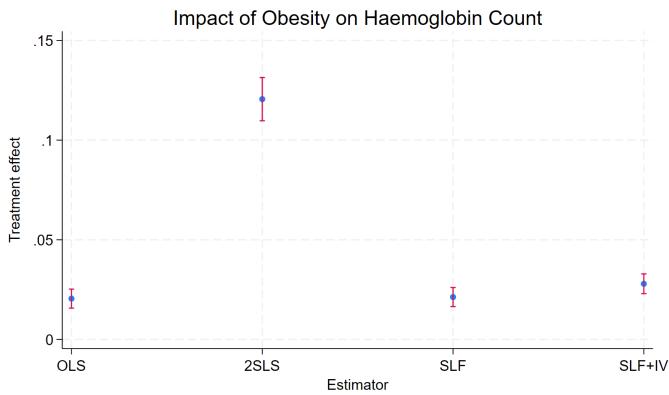


Fig 4: OLS, 2SLS, SLF, and SLF+IV estimates of the effect of obesity on haemoglobin level, with 95% confidence intervals

consumption of fish, chicken/meat consumption, and daily dried food consumption is significantly associated with a higher haemoglobin level. In Figure 4, the OLS, SLF, and SLF+IV are close to each other and have almost the same 95% confidence interval width. The 2SLS estimate has the largest magnitude and the widest confidence interval.

4.4.7 Virtue of shared latent factor models to control for omitted variables

The latent factor employed by both the shared latent factor models can be seen as an *index* of all the variables that either are unobservable or are observable but have been omitted from the analysis. The coefficient of the latent factor represents the combined effect of all such omitted variables on the treatment and outcome variables. So, if any essential control variable is deliberately dropped from the model, then the magnitude of the coefficient of the latent factor has to rise, as the omitted control variable becomes a part of the new latent factor.

We have run a new set of OLS, 2SLS, SLF, and SLF+IV models, wherein we deliberately dropped the control variable *place of residence: urban*. Table A3 in the Appendix summarizes the results of the new set of models. The magnitude of the coefficient of the latent factor has increased as the control variable has been dropped from the analysis for a given outcome variable and for both SLF

and SLF+IV. For instance, the magnitude of the coefficient of the latent factor for the outcome log of average diastolic reading rose from 21 (when place of residence: urban is included) to 118 (when place of residence: urban is excluded) under the SLF model, the magnitude of the coefficient of the latent factor for the outcome log of blood glucose level rose from 0.030 (when place of residence: urban is included) to 7.989 (when place of residence: urban is excluded) under the SLF+IV model. Hence, in this way, we have demonstrated how latent factor models capture the effects of all the omitted variables in order to render treatment effects with lower biases vis-a-vis OLS and 2SLS.

4.5 Discussion

In our empirical application, we find that obesity is significantly associated with higher diastolic and systolic readings, blood glucose levels, and haemoglobin levels under OLS, 2SLS, SLF, and SLF+IV models. This indicates that obese women are at a higher risk of being hypertensive and diabetic. On the other hand, obese women are at a lower risk of being anaemic. We find that the 2SLS estimates are very large for all the four outcomes. So, based on the theoretical claim of Banerjee and Basu (2021), the 2SLS estimates are more biased than the OLS, SLF, and SLF+IV estimates. We have demonstrated the virtue of the latent factor models in controlling for unobservable or omitted variables.

The coefficient of the latent factor is significant in both SLF and SLF+IV models, which indicates that there is unobserved heterogeneity between the treatment and control groups in both OLS and 2SLS models. We find that the magnitude of the coefficient of the latent factor is larger in the SLF model than in the SLF+IV model, which indicates that unobserved heterogeneity between the treatment and control groups is relatively a more serious concern in the OLS model than in the 2SLS model.

Our models have some limitations. Firstly, we could not control for the physical activities like exercising, yoga, etc., of the respondent due to the unavailability of information on such aspects in NFHS-5, although physical activities are an important determinant of obesity and blood pressure, blood sugar, and haemoglobin levels. Secondly, we did not control for the pregnancy and parity-related aspects of the women. Finally, we did not control whether the woman was already taking any medications for hypertension and/or diabetes. However, even if such variables have not been considered, the latent factor of the SLF and SLF+IV models shall capture these omitted variables to mitigate the bias of the estimates.

5 Conclusion

This paper highlights the choice between SLF and SLF+IV models based on IV strength and degree of contamination. Our theoretical and simulation results strongly suggest that the SLF+IV estimator is less (more) biased than the SLF estimator when the IV is strong (weak) and less (highly) contaminated. SLF and SLF+IV models perform better than naive OLS and 2SLS models, respectively, in terms of bias and MSE. We ascertained that SLF and SLF+IV models require a minimum of 4 and 2 outcome variables for parameter identification, respectively.

We have presented an empirical application wherein we analyzed the effect of obesity on the diastolic reading, systolic reading, blood glucose level and haemoglobin level, using OLS, 2SLS, SLF, and SLF+IV models and taking arm circumference as the IV. We found that all four models indicate that obesity is associated with a higher risk of hypertension and diabetes but a lower risk of anaemia. We demonstrated how the latent factor of the SLF and SLF+IV models capture the impact of variable(s) that are omitted in the traditional OLS and 2SLS models, respectively.

6 References

- Ahankari, A. S., Tata, L. J., & Fogarty, A. W. (2020). Weight, height, and midupper arm circumference are associated with haemoglobin levels in adolescent girls living in rural India: A cross-sectional study. *Maternal & Child Nutrition*, 16(2), e12908.
- Anderson, J. C., & Gerbing, D. W. (1988). Structural equation modeling in practice: A review and recommended two-step approach. *Psychological Bulletin*, 103(3), 411.
- Angrist, J. D., & Pischke, J. S. (2009). Mostly harmless econometrics: An empiricist's companion. Princeton University Press.
- Banerjee, S., & Basu, A. (2021). Estimating endogenous treatment effects using latent factor models with and without instrumental variables. *Econometrics*, 9(1), 14.
- Bollen, K. A., & Davis, W. R. (2009). Two rules of identification for structural equation models. *Structural Equation Modeling: A Multidisciplinary Journal*, 16(3), 523-536.
- Bollen, K., & Lennox, R. (1991). Conventional wisdom on measurement: A structural equation perspective. *Psychological Bulletin*, 110(2), 305.
- Bound, J., Jaeger, D., & Baker, R. (1993). The Cure Can Be Worse than the Disease: A Cautionary Tale Regarding Instrumental Variables. *National Bureau of Economic Research Technical Working Paper Series, No. 137*.
- Carneiro, P., Hansen, K. T., & Heckman J. J. (2001). Estimating Distributions Of Treatment Effects With An Application To The Returns To Schooling And Measurement Of The Effects Of Uncertainty On College Choice," *International Economic Review*, 2003, v44(2,May), 631-422.
- Cepeda-Lopez, A. C., & Baye, K. (2020). Obesity, iron deficiency and anaemia: a complex relationship. *Public Health Nutrition*, 23(10), 1703-1704.
- Craig, E., Bland, R., Ndirangu, J., & Reilly, J. J. (2014). Use of mid-upper arm circumference for determining overweight and overfatness in children and adolescents. *Archives of Disease in Childhood*, 99(8), 763-766.
- Crown, W. H., Henk, H. J., & Vanness, D. J. (2011). Some cautions on the use of instrumental variables estimators in outcomes research: how bias in instrumental variables estimators is affected by instrument strength, instrument contamination, and sample size. *Value in Health*, 14(8), 1078-1084.
- Danes, J. E., & Mann, O. K. (1984). Unidimensional measurement and structural equation models with latent variables. *Journal of Business Research*, 12(3), 337-352.
- Dereje, R., Girma, A., Molla, A., Simieneh, A. (2022). Mid upper arm circumference as screening tool of overweight or obesity among adult employees of Mizan Tepi University, Southwest Ethiopia. *Heliyon*, 8(10).
- Diamantopoulos, A., & Siguaw, J. A. (2006). Formative versus reflective indicators in organizational measure development: A comparison and empirical illustration. *British Journal of Management*, 17(4), 263-282.

- Diamantopoulos, A., & Winklhofer, H. M. (2001). Index construction with formative indicators: An alternative to scale development. *Journal of Marketing Research*, 38(2), 269-277.
- Forouhi, N. G., & Wareham, N. J. (2010). Epidemiology of diabetes. *Medicine*, 38(11), 602-606.
- Goldberger, A. S. (1972). Structural equation methods in the social sciences. *Econometrica: Journal of the Econometric Society*, 979-1001.
- Hauser, R. M., & Goldberger, A. S. (1971). The treatment of unobservable variables in path analysis. *Sociological Methodology*, 3, 81-117.
- International Institute for Population Sciences (IIPS) & ICF. (2021). National Family Health Survey (NFHS-5), 2019–21. http://rchiips.org/nfhs/NFHS-5Reports/NFHS-5_INDIA_REPORT.pdf
- Jordaan, E. M., Van den Berg, V. L., Van Rooyen, F. C., & Walsh, C. M. (2020). Obesity is associated with anaemia and iron deficiency indicators among women in the rural Free State, South Africa. *South African Journal of Clinical Nutrition*, 33(3), 72-78.
- Kiran, R., Harshitha., & Bhargava, M. (2022). Mid-upper arm circumference and neck circumference to screen for overweight-obesity in young adults in South India. *Helijon*, 8(12).
- Kumar, A., & Ram, F. (2013). Influence of family structure on child health: evidence from India. *Journal of Biosocial Science*, 45(5), 577-599.
- Land, K. C. (1969). Principles of path analysis. *Sociological Methodology*, 1, 3-37.
- Law, K. S., & Wong, C. S. (1999). Multidimensional constructs M structural equation analysis: An illustration using the job perception and job satisfaction constructs. *Journal of Management*, 25(2), 143-160.
- Lee, N., Cadogan, J. W., & Chamberlain, L. (2013). The MIMIC model and formative variables: problems and solutions. *AMS Review*, 3, 3-17.
- Lei, P. W., & Wu, Q. (2007). Introduction to structural equation modeling: Issues and practical considerations. *Educational Measurement: Issues and Practice*, 26(3), 33-43.
- Loenneke, J. P., Loprinzi, P. D., Abe, T., Thiebaud, R. S., Allen, K. M., Mouser, J. G., & Bemben, M. G. (2016). Arm circumference influences blood pressure even when applying the correct cuff size: Is a further correction needed?. *International Journal of Cardiology*, 202, 743-744.
- Maddala, G.S. (1992) Introduction to Econometrics. 2nd Edition, Prentice Hall, New Jersey.
- Mueller, R. O. (1997). Structural equation modeling: Back to basics. *Structural Equation Modeling: A Multidisciplinary Journal*, 4(4), 353-369.
- Nelson, C. R., & Startz, R. (1990). The Distribution of the Instrumental Variables Estimator and Its *t*-Ratio When the Instrument is a Poor One. *The Journal of Business*, 63(S1), S125.
- Posey, C., Roberts, T., Lowry, P. B., & Bennett, B. (2014). Multiple indicators and multiple causes (MIMIC) models as a mixed-modelling technique: A tutorial and an annotated example. *Communications of the Association for Information Systems*, 36(11).
- Saad, R. A., & Qutob, H. M. (2022). The relationship between anemia and obesity. *Expert Review of Hematology*, 15(10), 911-926.

- Segula, D. (2014). Complications of obesity in adults: a short review of the literature. *Malawi Medical Journal*, 26(1), 20-24.
- Staessen, J. A., Wang, J., Bianchi, G., & Birkenhäger, W. H. (2003). Essential hypertension. *The Lancet*, 361(9369), 1629-1641.
- Stein, C. M., Morris, N. J., & Nock, N. L. (2012). Structural equation modeling. *Statistical Human Genetics: Methods and Protocols*, 495-512.
- StataCorp. (2023). Stata Statistical Software: Release 18. College Station, TX: StataCorp LLC.
- Wang, J., He, L., Yang, N., Li, Z., Xu, L., Li, W., ... & Li, Y. (2022). Large mid-upper arm circumference is associated with reduced insulin resistance independent of BMI and waist circumference: A cross-sectional study in the Chinese population. *Frontiers in Endocrinology*, 13, 1054671.
- Wharton, S., Lau, D. C., Vallis, M., Sharma, A. M., Biertho, L., Campbell-Scherer, D., ... & Wicklum, S. (2020). Obesity in adults: a clinical practice guideline. *Cmaj*, 192(31), E875-E891.
- World Health Organization. 1998. Obesity: preventing and managing the global epidemic. Report of a WHO Consultation on Obesity. Geneva, Switzerland: World Health Organization.
- Wright, S. (1921). Correlation and causation. *Journal of Agricultural Research*, 20(7), 557.

7 Appendix

Relationship Between Strength of the IV and σ_z^2

The expression for the strength of an IV ($|Cov(D, z)|$) is given as follows:

$$\begin{aligned} |Cov(D, z)| &= |Cov(\pi z + \theta + u, z)| \\ &= |\pi\sigma_z^2 + Cov(z, \theta)| \end{aligned}$$

The expression for $|Cov(D, z)|$ shows that if $Cov(z, \theta)$ and π have the same sign, then the strength of the IV increases (decreases) as σ_z^2 increases (decreases), for given values of $Cov(z, \theta)$ and π . In other words, if values of $Cov(z, \theta)$ and π are given and have the same sign, then $|Cov(D, z)|$ becomes an *increasing* function of σ_z^2 .

Minimum number of outcome variables required when there are two IVs

Consider the following SLF+IV model with two IVs:

$$\begin{aligned} D &= \pi_1 z_1 + \pi_2 z_2 + \theta + u \\ Y^1 &= \beta_1^1 D + \gamma^1 \theta + w^1 \\ Y^2 &= \beta_1^2 D + \gamma^2 \theta + w^2 \\ &\vdots \\ Y^m &= \beta_1^m D + \gamma^m \theta + w^m \end{aligned}$$

We have made the following assumptions:

$$\begin{aligned}
& \text{Cov}(D, z_i) \neq 0, i = 1, 2 \\
& \text{Cov}(z_i, w^j) = 0, i = 1, 2, j = 1, 2, \dots, m \\
& \text{Cov}(z_i, u) = 0, i = 1, 2 \\
& \text{Cov}(\theta, u) = 0 \\
& \text{Cov}(\theta, w^i) = 0, i = 1, 2, \dots, m \\
& \text{Cov}(w^j, w^k) = 0, j \neq k, j = 1, 2, \dots, m, k = 1, 2, \dots, m \\
& \text{Cov}(u, w^i) = 0, i = 1, 2, \dots, m \\
& \text{Cov}(z_i, \theta) = 0, i = 1, 2
\end{aligned}$$

The covariance equations are as follows:

$$\begin{aligned}
\text{Cov}(Y^1, Y^2) &= \text{Cov}(\beta_1^1 D + \gamma^1 \theta + w^1, \beta_1^2 D + \gamma^2 \theta + w^2) \\
&= \beta_1^1 \beta_1^2 \sigma_D^2 + [\beta_1^1 \gamma^2 + \beta_1^2 \gamma^1 + \gamma^1 \gamma^2] \sigma_\theta^2 \\
\text{Cov}(Y^1, Y^3) &= \text{Cov}(\beta_1^1 D + \gamma^1 \theta + w^1, \beta_1^3 D + \gamma^3 \theta + w^3) \\
&= \beta_1^1 \beta_1^3 \sigma_D^2 + [\beta_1^1 \gamma^3 + \beta_1^3 \gamma^1 + \gamma^1 \gamma^3] \sigma_\theta^2 \\
&\vdots \\
\text{Cov}(Y^{m-1}, Y^m) &= \text{Cov}(\beta_1^{m-1} D + \gamma^{m-1} \theta + w^{m-1}, \beta_1^m D + \gamma^m \theta + w^m) \\
&= \beta_1^{m-1} \beta_1^m \sigma_D^2 + [\beta_1^{m-1} \gamma^m + \beta_1^m \gamma^{m-1} + \gamma^{m-1} \gamma^m] \sigma_\theta^2 \\
\text{Cov}(Y^1, D) &= \text{Cov}(\beta_1^1 D + \gamma^1 \theta + w^1, D) \\
&= \beta_1^1 \sigma_D^2 + \gamma^1 \sigma_\theta^2 \\
\text{Cov}(Y^2, D) &= \text{Cov}(\beta_1^2 D + \gamma^2 \theta + w^2, D) \\
&= \beta_1^2 \sigma_D^2 + \gamma^2 \sigma_\theta^2 \\
&\vdots \\
\text{Cov}(Y^m, D) &= \text{Cov}(\beta_1^m D + \gamma^m \theta + w^m, D) \\
&= \beta_1^m \sigma_D^2 + \gamma^m \sigma_\theta^2 \\
\text{Cov}(Y^1, z_1) &= \text{Cov}(\beta_1^1 D + \gamma^1 \theta + w^1, z_1) \\
&= \beta_1^1 \pi \sigma_{z_1}^2 \\
\text{Cov}(Y^2, z_1) &= \text{Cov}(\beta_1^2 D + \gamma^2 \theta + w^2, z_1) \\
&= \beta_1^2 \pi \sigma_{z_1}^2 \\
&\vdots \\
\text{Cov}(Y^m, z_1) &= \text{Cov}(\beta_1^m D + \gamma^m \theta + w^m, z_1) \\
&= \beta_1^m \pi \sigma_{z_1}^2
\end{aligned}$$

$$\begin{aligned}
Cov(Y^1, z_2) &= Cov(\beta_1^1 D + \gamma^1 \theta + w^1, z_2) \\
&= \beta_1^1 \pi \sigma_{z_2}^2 \\
Cov(Y^2, z_2) &= Cov(\beta_1^2 D + \gamma^2 \theta + w^2, z_2) \\
&= \beta_1^2 \pi \sigma_{z_2}^2 \\
&\vdots \\
Cov(Y^m, z_2) &= Cov(\beta_1^m D + \gamma^m \theta + w^m, z_2) \\
&= \beta_1^m \pi \sigma_{z_2}^2 \\
Cov(D, z_1) &= \pi \sigma_{z_1}^2 \\
Cov(D, z_2) &= \pi \sigma_{z_2}^2
\end{aligned}$$

There are $\frac{m(m-1)}{2}$ covariance equations for discordant pairs of outcomes, m covariance equations for each outcome variable and the treatment variable, $2m$ covariance equations for each IV and each outcome variable, and 2 covariance equations for each IV and the treatment variable. The parameters to be estimated include m “ β ” parameters, m “ γ ” parameters, 2 “ π ” parameters, and 1 “ σ_θ^2 ”. So, there are $\frac{m(m-1)}{2} + 3m + 2$ equations, and $2m + 3$ parameters to be estimated.

Number of covariance terms \geq Number of parameters to be estimated

$$\begin{aligned}
&\Rightarrow \frac{m(m-1)}{2} + 3m + 2 \geq 2m + 3 \\
&\Rightarrow m^2 + m - 2 \geq 0
\end{aligned}$$

Solving the above quadratic inequality, we get: $m \geq 1$. So, when there are two IVs, the SLF+IV model parameters can be estimated using only 1 outcome variable, vis-a-vis a minimum of 2 outcome variables required when there is only 1 IV.

Table A1(a): $Y^1, n = 1000$

		Y ¹		Degrees of IV contamination \rightarrow 0 (no contamination)								0.025		0.035		0.045		0.055		0.065		0.075	
True treatment effect = 2	IV Strength \downarrow	Baseline cases		Naive OLS		SLF		2SLS		SLF+IV		2SLS		SLF+IV		2SLS		SLF+IV		2SLS		SLF+IV	
std(z) = 0.16	Treatment effect estimate	3.760	2.420	1.798	2.001	5.189	2.232	5.831	2.482	6.433	4.278	6.980	6.055	7.280	6.304	7.568	7.636	7.505	7.505	202.505	202.505	-	
	First stage F-stat	-	-	15.659	-	51.006	-	72.541	-	99.203	-	128.535	-	161.127	-	215.18	-	278.39	-	221.82	-	-	
	% Bias	87.98	20.99	10.09	0.07	159.47	11.59	191.33	24.11	221.66	113.89	248.99	202.76	264.01	215.18	278.39	202.505	202.505	202.505	202.505	202.505	-	
	se	0.143	0.630	1.348	0.223	0.638	0.277	0.563	0.345	0.509	0.415	0.469	0.325	0.436	0.263	0.409	0.224	0.224	0.224	0.224	0.224	0.224	
	Coverage Probability	0.000	0.787	0.972	0.960	0.002	0.920	0.000	0.870	0.000	0.455	0.000	0.023	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	
	MSE	3.12	0.57	1.86	0.05	10.58	0.13	14.99	0.35	19.91	5.36	25.02	16.55	28.07	18.59	31.17	19.73	19.73	19.73	19.73	19.73	-	
std(z) = 0.2	Treatment effect estimate	3.760	2.420	1.871	2.000	4.398	2.209	5.013	2.325	5.487	2.522	5.912	3.438	6.423	3.535	6.763	6.010	6.010	6.010	6.010	6.010	-	
	First stage F-stat	-	-	29.962	-	67.967	10.47	150.66	-	119.898	-	149.502	-	172.489	-	211.13	167.49	238.16	238.16	200.50	200.50	-	
	% Bias	87.98	20.99	6.43	0.01	119.89	10.47	150.66	0.292	0.535	0.301	0.492	0.379	0.509	0.428	0.498	0.498	0.498	0.498	0.498	0.498	-	
	se	0.143	0.630	1.031	0.221	0.602	0.292	0.535	0.301	0.492	0.379	0.509	0.428	0.509	0.428	0.498	0.498	0.498	0.498	0.498	0.498	-	
	Coverage Probability	0.000	0.787	0.973	0.958	0.038	0.913	0.000	0.878	0.000	0.855	0.000	0.652	0.000	0.142	0.000	0.009	0.009	0.009	0.009	0.009	-	
	MSE	3.12	0.57	1.08	0.05	6.11	0.11	9.37	0.20	12.40	0.42	15.51	2.33	19.74	11.38	22.85	16.18	16.18	16.18	16.18	16.18	-	
std(z) = 0.24	Treatment effect estimate	3.760	2.420	1.909	2.000	3.832	2.197	4.384	2.297	4.844	2.428	5.218	2.611	5.536	3.009	5.926	4.334	4.334	4.334	4.334	4.334	-	
	First stage F-stat	-	-	43.938	-	94.975	-	115.598	-	136.902	-	157.298	-	183.114	-	-	-	-	-	-	-	-	-
	% Bias	87.98	20.99	4.56	0.00	91.58	9.85	119.21	14.85	142.21	21.39	160.90	30.53	176.80	50.43	196.32	116.71	116.71	116.71	116.71	116.71	-	
	se	0.143	0.630	0.844	0.218	0.566	0.252	0.511	0.279	0.471	0.329	0.442	0.434	0.418	0.549	0.398	0.578	0.578	0.578	0.578	0.578	-	
	Coverage Probability	0.000	0.787	0.970	0.959	0.118	0.909	0.005	0.855	0.000	0.825	0.000	0.817	0.000	0.767	0.000	0.394	0.394	0.394	0.394	0.394	-	
	MSE	3.12	0.57	0.72	0.05	3.68	0.10	5.95	0.17	8.31	0.29	10.55	0.56	12.68	1.32	15.57	5.78	5.78	5.78	5.78	5.78	-	
std(z) = 0.28	Treatment effect estimate	3.760	2.420	1.930	2.000	3.434	2.186	3.903	2.277	4.320	2.384	4.685	2.518	5.006	2.723	5.269	3.012	3.012	3.012	3.012	3.012	-	
	First stage F-stat	-	-	81.051	-	97.356	-	118.487	-	127.390	-	149.005	-	167.033	-	187.980	-	-	-	-	-	-	-
	% Bias	87.98	20.99	3.50	0.01	71.69	9.30	95.14	13.87	115.98	19.21	134.25	25.91	150.29	36.15	163.45	50.59	50.59	50.59	50.59	50.59	-	
	se	0.143	0.630	0.719	0.215	0.530	0.244	0.486	0.264	0.452	0.295	0.425	0.352	0.403	0.458	0.386	0.619	0.619	0.619	0.619	0.619	-	
	Coverage Probability	0.000	0.787	0.967	0.955	0.231	0.898	0.039	0.853	0.002	0.794	0.000	0.759	0.000	0.739	0.000	0.718	0.718	0.718	0.718	0.718	-	
	MSE	3.12	0.57	0.52	0.05	2.34	0.09	3.86	0.15	5.58	0.23	7.39	0.39	9.20	0.73	10.84	1.41	1.41	1.41	1.41	1.41	-	
std(z) = 0.32	Treatment effect estimate	3.760	2.420	1.945	2.000	3.149	2.176	3.546	2.259	3.910	2.352	4.241	2.465	4.546	2.723	5.269	3.012	3.012	3.012	3.012	3.012	-	
	First stage F-stat	-	-	118.443	-	166.738	-	183.950	-	205.485	-	232.178	-	233.649	-	256.335	-	-	-	-	-	-	-
	% Bias	87.98	20.99	2.77	0.02	57.44	8.81	77.32	12.94	95.49	17.61	112.03	23.24	127.28	36.32	140.86	39.94	39.94	39.94	39.94	39.94	-	
	se	0.143	0.630	0.629	0.212	0.496	0.236	0.461	0.251	0.433	0.273	0.410	0.307	0.391	0.367	0.375	0.469	0.469	0.469	0.469	0.469	-	
	Coverage Probability	0.000	0.787	0.966	0.960	0.369	0.898	0.103	0.853	0.012	0.779	0.000	0.728	0.000	0.687	0.000	0.649	0.649	0.649	0.649	0.649	-	
	MSE	3.12	0.57	0.40	0.04	1.57	0.09	2.60	0.13	3.84	0.20	5.19	0.31	6.63	0.50	8.08	0.86	0.86	0.86	0.86	0.86	-	
std(z) = 0.36	Treatment effect estimate	3.760	2.420	1.955	1.999	2.942	2.166	3.282	2.243	3.598	2.328	3.891	2.424	4.163	2.535	4.419	2.675	2.675	2.675	2.675	2.675	-	
	First stage F-stat	-	-	175.158	-	204.138	-	228.111	-	251.941	-	274.945	-	300.639	-	326.047	-	-	-	-	-	-	-
	% Bias	87.98	20.99	2.26	0.05	47.09	8.28	64.11	12.13	79.90	16.38	94.55	21.18	108.14	26.73	120.93	33.77	33.77	33.77	33.77	33.77	-	
	se	0.143	0.630	0.561	0.209	0.466	0.229	0.438	0.242	0.415	0.258	0.396	0.280	0.379	0.312	0.365	0.365	0.365	0.365	0.365	-		
	Coverage Probability	0.000	0.787	0.964	0.959	0.453	0.904	0.174	0.848	0.041	0.792	0.003	0.718	0.000	0.649	0.000	0.588	0.588	0.588	0.588	0.588	-	
	MSE	3.12	0.57	0.32	0.04	1.10	0.08	1.84	0.12	2.73	0.17	3.73	0.26	4.82	0.38	5.98	0.59	0.59	0.59	0.59	0.59	-	
std(z) = 1	Treatment effect estimate	3.760	2.420	1.988	1.997	2.169	2.075	2.241	2.107	2.312	2.139	2.382	2.171	2.452	2.205	2.522	2.238	2.238	2.238	2.238	2.238	-	
	First stage F-stat	-	-	749.168	-	797.245	-	818.697	-	836.864	-	851.653	-	871.746	-	-	-	-	-	-	-	-	
	% Bias	87.98	20.99	0.59	0.14	8.46	3.76	12.04	5.34	15.59	6.95	19.12	8.57	22.62	10.23	26.10	11.92	11.92	11.92	11.92	11.92	-	
	se	0.143	0.630	0.258	0.168	0.254	0.171	0.252	0.172	0.251	0.174	0.249	0.175	0.248	0.176	0.246	0.178	0.246	0.246	0.246	0.246	-	
	Coverage Probability	0.000	0.787	0.966	0.9510	0.929	0.905	0.842	0.908	0.752	0.867	0.830	0.545	0.779	0.440	0.779	0.440	0.779	0.440	0.779	0.440	-	
	MSE	3.12	0.57	0.07	0.03	0.09	0.03	0.12	0.04	0.16	0.05	0.21	0.06	0.27	0.07	0.33	0.09	0.33	0.09	0.33	0.09	-	

Notes:

First stage F-stat: F statistic for the first stage regression of 2SLS

% bias: percentage bias

se: standard error of treatment effect

MSE: Mean Squared Error

std(z): standard deviation of z

Table A1(b): $Y^2, n = 1000$

		True treatment effect = 3		Degrees of IV contamination \rightarrow 0 (no contamination)								0.025		0.035		0.045		0.055		0.065		0.075							
				Baseline cases		Naive OLS		SLF		SLF+IV		2SLS		SLF+IV		2SLS		SLF+IV		2SLS		SLF+IV							
std(z) = 0.16		Treatment effect estimate	4.762	3.423	2.790	3.005	6.183	3.236	6.805	3.490	7.436	5.278	7.980	7.054	8.280	7.303	8.565	7.435											
		First stage F-stat	-	15.659	14.11	-	51.006	-	72.541	16.32	147.85	99.203	-	128.535	-	161.127	-	202.505	-										
	% Bias	58.73	7.02	0.16	106.10	7.87	126.82	75.94	165.99	135.14	175.99	143.43	-	163.46	185.48	147.84	-	185.48	147.84										
	se	0.143	0.629	1.353	0.223	0.638	0.276	0.561	0.344	0.509	0.414	0.469	0.325	0.437	0.263	0.409	0.224	0.409	0.224										
	Coverage Probability	0.000	0.789	0.968	0.956	0.001	0.915	0.000	0.851	0.000	0.427	0.000	0.021	0.000	0.000	0.000	0.000	0.000	0.000	0.000									
	MSE	3.12	0.58	1.87	0.05	10.54	0.13	14.79	0.36	19.93	5.36	25.02	16.54	28.07	18.58	31.13	19.72												
std(z) = 0.2		Treatment effect estimate	4.762	3.423	2.867	3.005	5.394	3.213	6.008	3.327	6.477	3.525	6.901	4.442	7.430	6.352	7.754	7.009											
	First stage F-stat	-	29.962	14.11	4.44	0.18	79.81	7.12	100.28	10.91	119.898	17.51	130.04	48.06	147.66	111.75	158.46	132.64											
	% Bias	58.73	1.034	0.629	0.221	0.602	0.261	0.535	0.300	0.491	0.378	0.456	0.508	0.427	0.399	0.497	0.313												
	se	0.143	0.629	0.789	0.966	0.955	0.035	0.902	0.001	0.874	0.000	0.836	0.000	0.647	0.000	0.138	0.000	0.009	0.000										
	Coverage Probability	0.000	0.789	0.58	1.09	0.05	6.10	0.11	9.34	0.20	12.33	0.42	15.43	2.34	19.80	11.40	22.76	16.17											
	MSE	3.12	0.58	1.09	0.05	3.66	0.10	5.92	0.17	8.29	0.29	10.47	0.56	12.60	1.32	15.61	5.80												
std(z) = 0.24		Treatment effect estimate	4.762	3.423	2.905	3.004	4.829	3.200	5.380	3.301	5.840	3.431	6.206	3.613	6.524	4.013	6.930	5.337											
	First stage F-stat	-	43.938	3.16	0.13	60.95	6.68	65.68	5.66	79.33	10.03	94.66	14.38	106.85	20.44	117.48	33.75	131.01	77.92										
	% Bias	58.73	14.11	0.629	0.218	0.846	0.252	0.566	0.252	0.511	0.279	0.470	0.328	0.411	0.433	0.417	0.397	0.397	0.376										
	se	0.143	0.629	0.789	0.967	0.959	0.119	0.898	0.009	0.855	0.000	0.827	0.000	0.813	0.000	0.746	0.000	0.000	0.378										
	Coverage Probability	0.000	0.789	0.58	0.72	0.05	3.66	0.10	5.92	0.17	8.29	0.29	10.47	0.56	12.60	1.32	15.61	5.80											
	MSE	3.12	0.58	0.72	0.05	2.33	0.09	3.85	0.15	5.57	0.24	7.37	0.40	9.17	0.74	10.77	4.141												
std(z) = 0.28		Treatment effect estimate	4.762	3.423	2.927	3.002	4.430	3.189	4.900	3.280	5.317	3.387	5.681	3.524	6.001	3.728	6.259	4.013											
	First stage F-stat	-	81.051	2.43	0.08	47.68	6.29	63.34	9.34	77.23	12.90	89.36	17.38	100.04	24.26	108.62	33.77												
	% Bias	58.73	14.11	0.629	0.215	0.964	0.058	0.229	0.039	0.845	0.002	0.810	0.000	0.758	0.000	0.726	0.000	0.000	0.469										
	se	0.143	0.629	0.789	0.964	0.954	0.05	0.92	0.00	0.855	0.015	0.855	0.015	0.796	0.016	0.726	0.000	0.680	0.000	0.641									
	Coverage Probability	0.000	0.789	0.58	0.72	0.05	2.33	0.09	3.85	0.15	5.57	0.24	7.37	0.40	9.17	0.74	10.77	4.141											
	MSE	3.12	0.58	0.72	0.05	2.33	0.09	3.85	0.15	5.57	0.24	7.37	0.40	9.17	0.74	10.77	4.141												
std(z) = 0.32		Treatment effect estimate	4.762	3.423	2.943	3.003	4.147	3.179	4.544	3.262	4.907	3.356	5.238	3.468	5.542	3.609	5.815	3.801											
	First stage F-stat	-	118.443	1.91	0.11	38.23	5.98	51.47	8.72	63.58	11.85	74.61	15.61	84.75	20.31	93.83	26.335	33.77											
	% Bias	58.73	14.11	0.629	0.212	0.964	0.051	0.355	0.900	0.101	0.852	0.016	0.796	0.020	0.726	0.000	0.680	0.000	0.641										
	se	0.143	0.629	0.789	0.964	0.954	0.04	0.92	0.00	0.855	0.015	0.855	0.015	0.796	0.020	0.726	0.000	0.680	0.000	0.641									
	Coverage Probability	0.000	0.789	0.58	0.72	0.04	1.56	0.56	2.60	0.13	3.83	0.20	5.18	0.31	6.62	0.51	8.06	0.86											
	MSE	3.12	0.58	0.72	0.04	1.56	0.56	2.60	0.13	3.83	0.20	5.18	0.31	6.62	0.51	8.06	0.86												
std(z) = 0.36		Treatment effect estimate	4.762	3.423	2.954	3.003	3.940	3.170	4.280	3.246	4.596	3.331	4.889	3.427	5.161	3.538	5.417	3.679											
	First stage F-stat	-	175.158	1.54	0.10	31.33	5.65	42.67	8.21	53.21	11.04	62.96	14.25	72.02	17.95	80.55	22.64												
	% Bias	58.73	14.11	0.629	0.209	0.963	0.947	0.466	0.229	0.438	0.241	0.415	0.288	0.396	0.280	0.379	0.312	0.364											
	se	0.143	0.629	0.789	0.963	0.957	0.32	0.40	0.10	0.83	0.12	0.72	0.18	3.72	0.26	0.712	0.001	0.642	0.000	0.573									
	Coverage Probability	0.000	0.789	0.58	0.72	0.03	0.97	0.09	0.08	0.12	0.04	0.12	0.04	0.160	0.051	0.21	0.06	0.27	0.07	0.33	0.09								
	MSE	3.12	0.58	0.72	0.03	0.97	0.09	0.08	0.12	0.04	0.12	0.04	0.160	0.051	0.21	0.06	0.27	0.07	0.33	0.09									
std(z) = 1		Treatment effect estimate	4.762	3.423	2.988	3.001	3.169	3.079	3.240	3.111	3.311	3.143	3.382	3.175	3.452	3.208	3.522	3.242											
	First stage F-stat	-	749.168	0.40	0.03	5.63	2.62	8.02	8.02	8.69	8.64	836.864	-	851.653	-	871.746	-	876.356	-										
	% Bias	58.73	14.11	0.629	0.258	0.168	0.254	0.171	0.252	0.172	0.251	0.174	0.249	0.175	0.248	0.177	0.246	0.178	0.246	0.178									
	se	0.143	0.629	0.789	0.957	0.9480	0.904	0.927	0.838	0.907	0.758	0.872	0.666	0.836	0.557	0.798	0.441	0.721	0.59										
	Coverage Probability	0.000	0.789	0.58	0.72	0.03	0.97	0.09	0.08	0.12	0.04	0.12	0.04	0.160	0.051	0.21	0.06	0.27	0.07	0.33	0.09								
	MSE	3.12	0.58	0.72	0.03	0.97	0.09	0.08	0.12	0.04	0.12	0.04	0.160	0.051	0.21	0.06	0.27	0.07	0.33	0.09									

Notes:

First stage F-stat: F statistic for the first stage regression of 2SLS

% bias: percentage bias

se: standard error of treatment effect

MSE: Mean Squared Error

std(z): standard deviation of z

Table A1(c): $Y^3, n=1000$

		Y^3		Degrees of IV contamination $\rightarrow 0$ (no contamination)								0.025		0.035		0.045		0.055		0.065		0.075		
True treatment effect = 1 IV Strength \downarrow		Baseline cases		Naive OLS		SLF		2SLS		SLF+IV		2SLS		SLF+IV		2SLS		SLF+IV		2SLS		SLF+IV		
std(z) = 0.16	Treatment effect estimate	2.761	1.419	0.771	1.008	4.179	1.239	4.813	1.495	5.447	3.284	5.973	5.055	6.274	5.303	6.562	5.437	-	-	202.505	-	-	-	
	First stage F-stat	-	-	15.659	-	51.006	-	72.541	-	99.203	-	128.535	-	161.127	-	161.127	-	202.505	-	-	-	-	-	
	% Bias	176.06	41.92	0.83	317.85	23.94	49.53	444.74	228.44	0.561	0.344	0.414	0.468	0.324	0.436	0.262	0.409	0.224	0.000	0.000	0.000	0.000	0.000	
	se	0.143	0.631	1.357	0.223	0.637	0.276	0.561	0.309	0.000	0.000	0.000	0.000	0.023	0.023	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	
	Coverage Probability	0.000	0.785	0.972	0.57	1.89	0.05	10.51	0.13	14.86	0.36	20.04	5.39	24.95	16.55	28.01	18.59	31.10	19.74	-	-	-	-	-
	MSE	3.12	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-
std(z) = 0.2	Treatment effect estimate	2.761	1.419	0.854	1.008	3.388	1.216	4.005	1.332	4.474	1.530	4.914	2.445	5.422	4.351	5.753	5.008	-	-	-	-	-	-	-
	First stage F-stat	-	-	29.962	-	67.967	-	95.004	-	119.898	-	149.502	-	172.489	-	210.510	-	475.27	-	400.82	-	-	-	-
	% Bias	176.06	41.92	14.64	0.81	238.76	21.61	300.53	33.20	347.45	52.99	391.39	144.51	442.21	335.08	-	-	-	-	-	-	-	-	-
	se	0.143	0.631	1.035	0.221	0.601	0.261	0.534	0.300	0.491	0.378	0.456	0.509	0.427	0.399	0.406	0.406	-	-	-	-	-	-	-
	Coverage Probability	0.000	0.785	0.969	0.962	0.027	0.901	0.000	0.868	0.000	0.828	0.000	0.634	0.000	0.140	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	
	MSE	3.12	0.57	1.09	0.05	6.06	0.11	9.32	0.20	12.31	0.42	15.53	2.35	19.74	11.39	22.75	16.16	-	-	-	-	-	-	-
std(z) = 0.24	Treatment effect estimate	2.761	1.419	0.896	1.007	2.822	1.203	3.376	1.305	3.837	1.437	4.206	1.618	4.526	2.020	4.934	3.338	-	-	-	-	-	-	-
	First stage F-stat	-	-	43.938	-	94.975	-	115.598	-	136.902	-	157.298	-	183.114	-	-	-	-	-	-	-	-	-	-
	% Bias	176.06	41.92	10.41	0.70	182.20	20.31	237.37	30.47	283.70	43.65	320.56	61.84	322.58	102.02	393.39	102.02	-	-	-	-	-	-	-
	se	0.143	0.631	0.846	0.218	0.566	0.252	0.511	0.279	0.470	0.328	0.441	0.418	0.434	0.548	0.398	0.577	0.392	0.392	0.392	0.392	0.392	0.392	
	Coverage Probability	0.000	0.785	0.968	0.965	0.122	0.289	0.007	0.864	0.000	0.813	0.000	0.793	0.000	0.755	0.000	0.744	0.000	0.744	0.000	0.744	0.000	0.744	
	MSE	3.12	0.57	0.73	0.05	3.64	0.10	5.90	0.17	8.27	0.30	10.47	0.57	12.61	1.34	15.63	10.63	15.63	5.80	-	-	-	-	
std(z) = 0.28	Treatment effect estimate	2.761	1.419	0.920	1.007	2.425	1.193	2.895	1.285	3.312	1.332	3.680	1.526	3.997	1.731	-	-	-	-	-	-	-	-	-
	First stage F-stat	-	-	81.051	-	97.356	-	118.487	-	127.390	-	149.005	-	167.033	-	187.980	-	-	-	-	-	-	-	-
	% Bias	176.06	41.92	8.02	0.71	142.49	19.31	189.53	28.48	231.25	39.17	268.04	52.62	299.74	73.11	325.70	73.11	101.91	101.91	-	-	-	-	-
	se	0.143	0.631	0.720	0.215	0.530	0.243	0.485	0.263	0.452	0.295	0.425	0.351	0.403	0.458	0.386	0.458	0.386	0.617	0.617	-	-	-	
	Coverage Probability	0.000	0.785	0.968	0.960	0.231	0.389	0.034	0.849	0.004	0.758	0.000	0.752	0.000	0.752	0.000	0.744	0.000	0.744	0.000	0.744	0.000	0.744	
	MSE	3.12	0.57	0.52	0.05	0.52	0.21	0.31	0.10	0.383	0.15	0.55	0.24	0.736	0.40	0.915	0.74	10.76	1.42	-	-	-	-	
std(z) = 0.32	Treatment effect estimate	2.761	1.419	0.937	1.006	2.143	1.183	2.540	1.266	2.904	1.360	3.235	1.472	3.540	1.613	3.812	1.806	-	-	-	-	-	-	-
	First stage F-stat	-	-	118.443	-	166.738	-	183.950	-	205.485	-	232.178	-	233.649	-	256.335	-	-	-	-	-	-	-	-
	% Bias	176.06	41.92	6.34	0.62	114.28	18.34	154.03	26.59	190.39	35.99	223.48	47.22	244.03	61.35	281.22	61.35	80.56	80.56	-	-	-	-	-
	se	0.143	0.631	0.630	0.211	0.496	0.236	0.461	0.251	0.433	0.272	0.410	0.306	0.391	0.367	0.374	0.374	0.469	-	-	-	-	-	
	Coverage Probability	0.000	0.785	0.967	0.953	0.339	0.900	0.099	0.843	0.012	0.775	0.001	0.717	0.000	0.672	0.000	0.672	0.000	0.635	0.000	0.635	0.000		
	MSE	3.12	0.57	0.40	0.04	1.55	0.09	2.59	0.13	3.81	0.20	5.16	0.32	6.61	0.51	8.05	0.87	-	-	-	-	-	-	-
std(z) = 0.36	Treatment effect estimate	2.761	1.419	0.949	1.007	1.936	1.174	2.277	1.251	2.593	1.336	2.886	1.431	3.158	1.543	3.414	1.684	-	-	-	-	-	-	-
	First stage F-stat	-	-	175.158	-	204.138	-	228.111	-	251.941	-	274.945	-	300.699	-	312.047	-	-	-	-	-	-	-	-
	% Bias	176.06	41.92	5.12	0.69	93.62	17.38	127.67	25.09	159.30	33.59	188.60	43.14	215.81	54.31	241.38	54.31	68.39	68.39	-	-	-	-	
	se	0.143	0.631	0.562	0.208	0.466	0.229	0.438	0.241	0.415	0.288	0.395	0.258	0.379	0.312	0.364	0.364	0.364	-	-	-	-	-	
	Coverage Probability	0.000	0.785	0.968	0.958	0.449	0.904	0.185	0.838	0.036	0.762	0.006	0.700	0.000	0.629	0.000	0.629	0.000	0.560	0.000	0.560	0.000		
	MSE	3.12	0.57	0.32	0.04	1.09	0.08	1.82	0.12	2.71	0.18	3.71	0.26	4.80	0.39	5.96	0.39	5.96	0.39	-	-	-	-	
std(z) = 1	Treatment effect estimate	2.761	1.419	0.986	0.999	1.167	1.077	1.239	1.109	1.310	1.141	1.380	1.173	1.450	1.206	1.520	1.240	-	-	-	-	-	-	-
	First stage F-stat	-	-	749.168	-	797.245	-	818.697	-	836.864	-	851.653	-	871.746	-	876.356	-	-	-	-	-	-	-	-
	% Bias	176.06	41.92	1.37	0.09	7.66	7.66	23.88	10.85	30.98	14.05	38.03	17.28	45.04	20.60	51.99	23.98	-	-	-	-	-	-	-
	se	0.143	0.631	0.258	0.168	0.254	0.171	0.252	0.251	0.251	0.174	0.175	0.248	0.177	0.246	0.177	0.246	0.177	-	-	-	-	-	-
	Coverage Probability	0.000	0.785	0.970	0.9450	0.892	0.933	0.847	0.914	0.776	0.870	0.674	0.831	0.547	0.782	0.426	0.782	0.426	-	-	-	-	-	-
	MSE	3.12	0.57	0.07	0.03	0.09	0.04	0.12	0.04	0.16	0.05	0.21	0.06	0.26	0.07	0.33	0.09	0.33	0.09	-	-	-	-	

Notes:

First stage F-stat: F statistic for the first stage regression of 2SLS

% bias: percentage bias

se: standard error of treatment effect

MSE: Mean Squared Error

std(z): standard deviation of z

Table A(d): $Y^4, n = 1000$

		True treatment effect = 1		Degrees of IV contamination \rightarrow		0 (no contamination)		0.025		0.035		0.045		0.055		0.065		0.075				
		IV Strength \downarrow		Baseline cases		Naive OLS		SLF		SLF+IV		2SLS		SLF+IV		2SLS		SLF+IV		2SLS		
std(z) = 0.16		Treatment effect estimate	2.756	1.413	0.785	4.178	1.234	4.803	1.488	5.415	3.275	5.973	5.052	6.275	5.300	6.560	5.432					
		First stage F-stat			15.659	51.006	72.541	99.203	128.535	161.127	497.29	405.24	527.54	429.97	555.98	429.97	202.505					
	% Bias		175.58	41.35	21.46	0.23	317.77	23.40	48.82	441.48	227.47	497.29	405.24	527.54	429.97	555.98	429.97					
	se		0.143	0.631	1.350	0.223	0.638	0.276	0.561	0.344	0.509	0.414	0.468	0.325	0.436	0.263	0.409	0.224				
	Coverage Probability		0.000	0.697	0.974	0.951	0.002	0.906	0.000	0.857	0.000	0.440	0.000	0.023	0.000	0.000	0.000	0.000	0.000	0.000	0.000	
	MSE		3.10	0.57	1.87	0.05	10.50	0.13	14.78	0.36	19.75	5.35	24.95	16.53	28.02	18.56	31.08	19.70				
std(z) = 0.2		Treatment effect estimate	2.756	1.413	0.863	1.002	3.391	1.211	4.005	1.326	4.472	1.523	4.905	2.437	5.416	4.350	5.751	5.006				
		First stage F-stat			175.58	41.35	13.70	0.15	239.11	21.06	300.49	32.56	347.22	52.26	390.55	143.69	441.57	335.00	210.510			
	% Bias		0.143	0.631	1.032	0.221	0.602	0.261	0.535	0.301	0.491	0.378	0.456	0.509	0.427	0.399	0.407	0.313				
	se		0.000	0.697	0.973	0.954	0.044	0.901	0.000	0.867	0.000	0.842	0.000	0.645	0.000	0.148	0.000	0.009				
	Coverage Probability		0.000	0.57	1.08	0.05	6.08	0.11	9.32	0.20	12.30	0.42	15.46	2.32	19.68	11.38	22.74	16.15				
	MSE		3.10	0.57	1.08	0.05	3.65	0.10	5.91	0.17	8.26	0.29	10.48	0.56	12.56	1.32	15.52	10.73				
std(z) = 0.24		Treatment effect estimate	2.756	1.413	0.901	1.001	2.824	1.197	3.377	1.297	3.836	1.429	4.206	1.609	4.519	2.011	4.920	3.336				
		First stage F-stat			43.938	5	9.92	0.07	182.43	19.67	237.73	29.74	283.58	42.87	320.64	60.87	351.95	101.06				
	% Bias		175.58	41.35	0.631	0.845	0.218	0.566	0.252	0.511	0.279	0.471	0.328	0.441	0.418	0.343	0.344	0.398				
	se		0.143	0.631	0.697	0.972	0.957	0.135	0.898	0.008	0.859	0.000	0.830	0.000	0.817	0.000	0.761	0.000	0.386			
	Coverage Probability		0.000	0.57	0.72	0.05	3.65	0.10	5.91	0.17	8.26	0.29	10.48	0.56	12.56	1.32	15.52	10.73				
	MSE		3.10	0.57	0.72	0.05	2.31	0.09	3.83	0.15	5.55	0.23	7.35	0.39	9.13	0.73	10.73	1.40				
std(z) = 0.28		Treatment effect estimate	2.756	1.413	0.923	1.000	2.426	1.185	2.895	1.276	3.312	1.383	3.678	1.519	3.995	1.724	4.252	2.009				
		First stage F-stat			81.051	5	97.356	5	118.487	5	127.390	5	149.005	5	167.033	5	187.980					
	% Bias		175.58	41.35	7.68	0.01	142.57	18.52	189.53	27.60	231.20	38.30	267.78	51.88	299.51	60.64	325.21	100.89				
	se		0.143	0.631	0.720	0.215	0.530	0.243	0.486	0.263	0.452	0.295	0.425	0.352	0.404	0.458	0.386	0.617				
	Coverage Probability		0.000	0.697	0.972	0.953	0.240	0.892	0.046	0.857	0.001	0.806	0.000	0.775	0.000	0.737	0.000	0.719				
	MSE		3.10	0.57	0.52	0.05	2.31	0.09	3.83	0.15	5.55	0.23	7.35	0.39	9.13	0.73	10.73	1.40				
std(z) = 0.32		Treatment effect estimate	2.756	1.413	0.939	0.999	2.144	1.176	2.541	1.258	2.904	1.352	3.234	1.465	3.539	1.606	3.811	1.799				
		First stage F-stat			118.443	6.11	0.08	114.39	17.55	183.950	5	205.485	5	232.178	5	233.649	5	236.335				
	% Bias		175.58	41.35	0.631	0.630	0.211	0.496	0.236	0.461	0.251	0.433	0.273	0.410	0.306	0.391	0.367	0.375				
	se		0.143	0.631	0.697	0.973	0.950	0.345	1.55	0.09	2.59	0.112	0.846	0.014	0.802	0.000	0.741	0.000	0.690	0.000	0.659	
	Coverage Probability		0.000	0.57	0.57	0.04	0.40	0.04	1.55	0.09	3.81	0.20	5.16	0.31	6.60	0.50	8.04	0.86				
	MSE		3.10	0.57	0.57	0.05	2.31	0.09	3.83	0.15	5.55	0.23	7.35	0.39	9.13	0.73	10.73	1.40				
std(z) = 0.36		Treatment effect estimate	2.756	1.413	0.950	0.999	1.937	1.165	2.277	1.242	2.593	1.327	2.885	1.422	3.157	1.534	3.413	1.675				
		First stage F-stat			175.58	41.35	4.98	0.15	93.66	16.46	127.69	24.22	159.28	32.69	188.55	42.24	215.72	53.38	241.32			
	% Bias		0.143	0.631	0.697	0.971	0.951	0.440	0.895	0.191	0.852	0.050	0.796	0.004	0.729	0.000	0.654	0.000	0.596			
	se		0.000	0.697	0.971	0.57	0.32	0.04	1.09	0.08	1.82	0.12	2.71	0.17	3.71	0.26	4.80	0.38	5.96	0.59		
	Coverage Probability		0.000	0.697	0.960	0.944	0.905	0.928	0.849	0.905	0.764	0.871	0.650	0.847	0.547	0.804	0.422	0.743				
	MSE		3.10	0.57	0.57	0.07	0.03	0.09	0.03	0.12	0.04	0.16	0.05	0.21	0.06	0.07	0.03	0.33	0.09			

Notes:

- First stage F-stat: F statistic for the first stage regression of 2SLS
- % bias: percentage bias
- se: standard error of treatment effect
- MSE: Mean Squared Error
- std(z): standard deviation of z

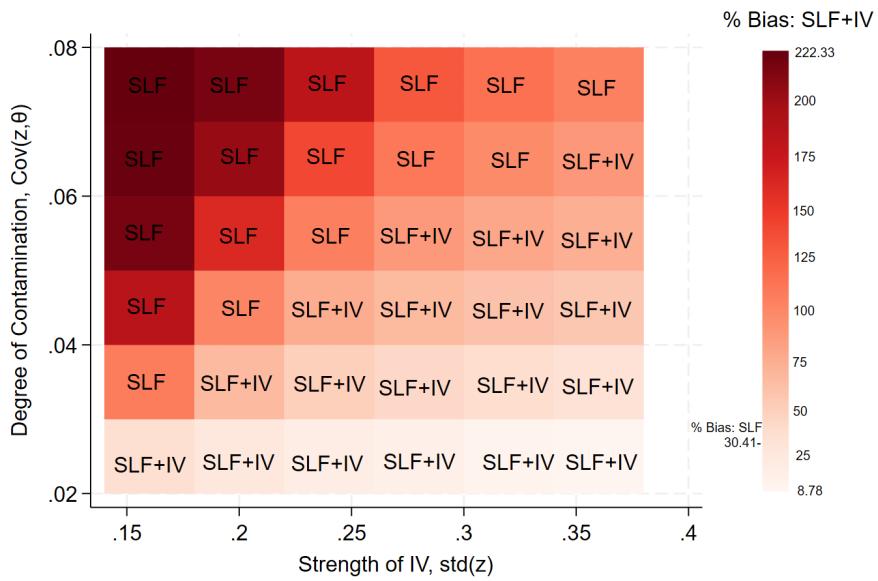


Fig A1: Heatplot for percentage bias of SLF+IV estimator for Y^1 , $n = 500$; different combinations of $Cov(z, \theta)$ and $std(z)$; colour intensity depicts magnitude of the bias

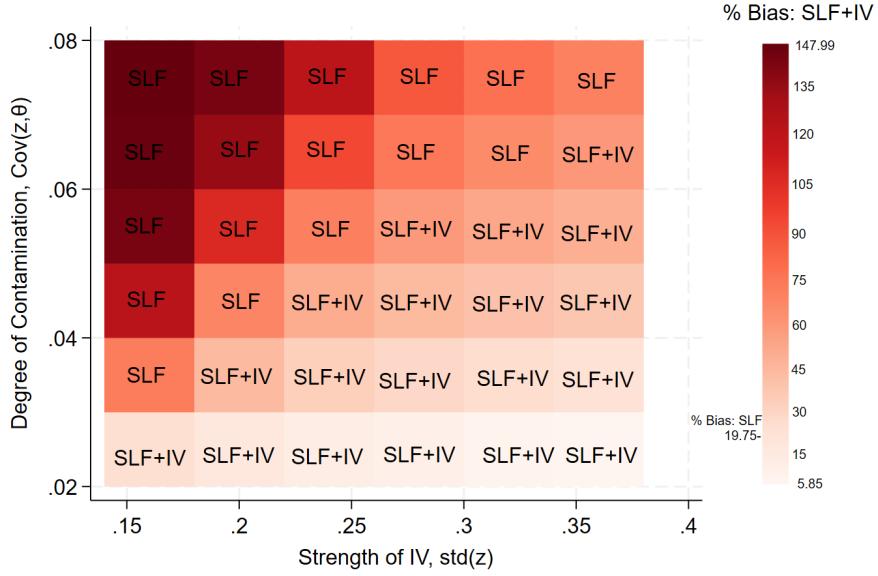


Fig A2: Heatplot for percentage bias of SLF+IV estimator for Y^2 , $n = 500$; different combinations of $Cov(z, \theta)$ and $std(z)$; colour intensity depicts magnitude of the bias

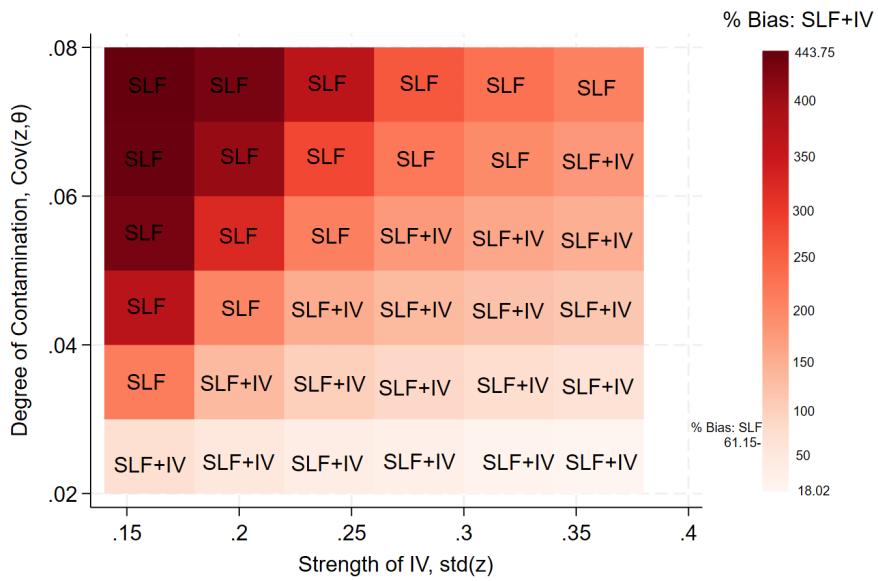


Fig A3: Heatplot for percentage bias of SLF+IV estimator for Y^3 , $n = 500$; different combinations of $Cov(z, \theta)$ and $std(z)$; colour intensity depicts magnitude of the bias

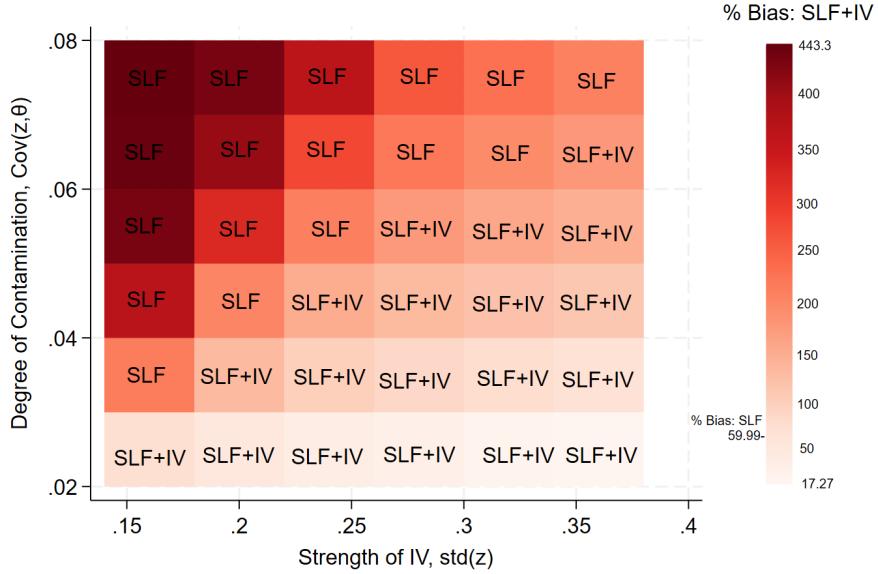


Fig A4: Heatplot for percentage bias of SLF+IV estimator for Y^4 , $n = 500$; different combinations of $Cov(z, \theta)$ and $std(z)$; colour intensity depicts magnitude of the bias

Table A2: Effect of obesity on diastolic reading, systolic reading, blood glucose level, and haemoglobin level, under OLS, 2SLS, SLE, and SLE+IV models

VARIABLES	log of average diastolic reading				log of average systolic reading				log of blood glucose level				log of haemoglobin level			
	OLS	2SLS	SLF	SLF+IV	OLS	2SLS	SLF	SLF+IV	OLS	2SLS	SLF	SLF+IV	OLS	2SLS	SLF	SLF+IV
Treatment variable: obese	0.0570*** (0.00191)	0.1105*** (0.00398)	0.0577*** (0.00194)	0.0841*** (0.00320)	0.0272*** (0.00199)	0.0492*** (0.00145)	0.0272*** (0.00139)	0.0505*** (0.00038)	0.0785*** (0.00037)	0.0978*** (0.00013)	0.0691*** (0.00018)	0.0691*** (0.00018)	0.1212*** (0.00043)	0.0295*** (0.00043)	0.1212*** (0.00043)	0.0295*** (0.00043)
Place of residence: urban	-0.006531*** (0.00135)	-0.006510*** (0.00139)	-0.006510*** (0.00138)	-0.006510*** (0.00135)	-0.007814*** (0.00140)	-0.007814*** (0.00140)	-0.007814*** (0.00136)	-0.006965*** (0.00035)	-0.006965*** (0.00022)	-0.006965*** (0.00022)	-0.006965*** (0.00022)	-0.006965*** (0.00022)	0.0114*** (0.00181)	0.0114*** (0.00181)	0.0114*** (0.00181)	0.0114*** (0.00181)
Current Age	0.002917*** (0.000202)	0.003018*** (0.000203)	0.003018*** (0.000202)	0.003018*** (0.000202)	0.003265*** (0.000203)	0.003288*** (0.000203)	0.003288*** (0.000203)	0.003288*** (0.000203)	0.003579*** (0.000203)	0.003579*** (0.000203)	0.003579*** (0.000203)	0.003579*** (0.000203)	0.0114*** (0.00181)	0.0114*** (0.00181)	0.0114*** (0.00181)	0.0114*** (0.00181)
Alcohol consumption: yes	0.0271*** (0.00028)	0.0275*** (0.00028)	0.0275*** (0.00027)	0.0275*** (0.00027)	0.0274*** (0.00027)	0.0368*** (0.00027)	0.0368*** (0.00027)	0.0373*** (0.00027)	0.0373*** (0.00027)	0.0373*** (0.00027)	0.0373*** (0.00027)	0.0373*** (0.00027)	0.0672*** (0.00019)	0.0672*** (0.00019)	0.0672*** (0.00019)	0.0672*** (0.00019)
Milk/Curd consumption: daily	0.000398 (0.00242)	-0.000858 (0.00240)	0.000398 (0.00237)	0.000398 (0.00234)	0.00725*** (0.00023)	0.00801*** (0.00023)	0.00719*** (0.00023)	0.00719*** (0.00023)	0.00697*** (0.00023)	0.00697*** (0.00023)	0.00697*** (0.00023)	0.00697*** (0.00023)	0.00808*** (0.00036)	0.00808*** (0.00036)	0.00808*** (0.00036)	0.00808*** (0.00036)
Milk/Curd consumption: weekly	0.000778 (0.00219)	0.001139 (0.00250)	0.000777 (0.00249)	0.001138 (0.00249)	0.0120*** (0.00024)	0.0109*** (0.00024)	0.0120*** (0.00024)	0.0109*** (0.00024)	0.00968*** (0.00024)	0.00968*** (0.00024)	0.00968*** (0.00024)	0.00968*** (0.00024)	0.00965*** (0.00036)	0.00965*** (0.00036)	0.00965*** (0.00036)	0.00965*** (0.00036)
Milk/Curd consumption: occasionally	0.000187 (0.00246)	0.000188 (0.00246)	0.000187 (0.00246)	0.000187 (0.00246)	-0.2962*** (0.00026)	-1.70e-05 (0.00026)	-0.005242*** (0.00026)	-0.005242*** (0.00026)	0.00506*** (0.00026)	0.00506*** (0.00026)	0.00506*** (0.00026)	0.00506*** (0.00026)	-0.00968*** (0.00037)	-0.00968*** (0.00037)	-0.00968*** (0.00037)	-0.00968*** (0.00037)
Pulses/Beans consumption: daily	0.0259*** (0.0121)	0.0259*** (0.0121)	0.0259*** (0.0121)	0.0259*** (0.0121)	0.0259*** (0.0121)	0.0259*** (0.0121)	0.0259*** (0.0121)	0.0259*** (0.0121)	0.0259*** (0.0121)	0.0259*** (0.0121)	0.0259*** (0.0121)	0.0259*** (0.0121)	0.0259*** (0.0121)	0.0259*** (0.0121)	0.0259*** (0.0121)	0.0259*** (0.0121)
Pulses/Beans consumption: weekly	0.0294** (0.0125)	0.0294** (0.0125)	0.0294** (0.0125)	0.0294** (0.0125)	0.0294** (0.0125)	0.0294** (0.0125)	0.0294** (0.0125)	0.0294** (0.0125)	0.0294** (0.0125)	0.0294** (0.0125)	0.0294** (0.0125)	0.0294** (0.0125)	0.0294** (0.0125)	0.0294** (0.0125)	0.0294** (0.0125)	0.0294** (0.0125)
Pulses/Beans consumption: occasionally	0.0322** (0.0126)	0.0322** (0.0126)	0.0322** (0.0126)	0.0322** (0.0126)	0.0328*** (0.0126)	0.0263*** (0.0126)	0.0263*** (0.0126)	0.0263*** (0.0126)	0.0263*** (0.0126)	0.0263*** (0.0126)	0.0263*** (0.0126)	0.0263*** (0.0126)	0.0225*** (0.0130)	0.0225*** (0.0130)	0.0225*** (0.0130)	0.0225*** (0.0130)
Dark Green Leafy Vegetables consumption: daily	0.00544 (0.0204)	0.00544 (0.0204)	0.00544 (0.0204)	0.00544 (0.0204)	0.010428 (0.00026)	-0.0112 (0.00026)	-0.0112 (0.00026)	-0.0112 (0.00026)	-0.0122 (0.00026)	-0.0122 (0.00026)	-0.0122 (0.00026)	-0.0122 (0.00026)	-0.0148e-05 (0.00026)	-0.0148e-05 (0.00026)	-0.0148e-05 (0.00026)	-0.0148e-05 (0.00026)
Dark Green Leafy Vegetables consumption: weekly	0.00196 (0.0204)	-0.000650 (0.0204)	0.000556 (0.0204)	-0.000908 (0.0204)	-0.0176 (0.00026)	-0.0186 (0.00026)	-0.0176 (0.00026)	-0.0176 (0.00026)	-0.0212 (0.00026)	-0.0212 (0.00026)	-0.0212 (0.00026)	-0.0212 (0.00026)	-0.00608 (0.00026)	-0.00608 (0.00026)	-0.00608 (0.00026)	-0.00608 (0.00026)
Dark Green Leafy Vegetables consumption: occasionally	-0.00353 (0.0205)	-0.002829 (0.0205)	-0.002829 (0.0205)	-0.002829 (0.0205)	-0.00848 (0.00026)	-0.01221 (0.00026)	-0.01221 (0.00026)	-0.01221 (0.00026)	-0.0231 (0.00026)	-0.0231 (0.00026)	-0.0231 (0.00026)	-0.0231 (0.00026)	-0.01619 (0.00026)	-0.01619 (0.00026)	-0.01619 (0.00026)	-0.01619 (0.00026)
Fruits consumption: daily	0.000591 (0.00355)	-0.010404 (0.00355)	-0.010404 (0.00355)	-0.010404 (0.00355)	-0.006766 (0.00026)	-0.006666 (0.00026)	-0.006666 (0.00026)	-0.006666 (0.00026)	-0.007579 (0.00026)	-0.007579 (0.00026)	-0.007579 (0.00026)	-0.007579 (0.00026)	-0.00876 (0.00026)	-0.00876 (0.00026)	-0.00876 (0.00026)	-0.00876 (0.00026)
Fruits consumption: weekly	0.00296 (0.00256)	0.00296 (0.00256)	0.00296 (0.00256)	0.00296 (0.00256)	0.00297 (0.00256)	0.00297 (0.00256)	0.00297 (0.00256)	0.00297 (0.00256)	0.00297 (0.00256)	0.00297 (0.00256)	0.00297 (0.00256)	0.00297 (0.00256)	0.02026 (0.00026)	0.02026 (0.00026)	0.02026 (0.00026)	0.02026 (0.00026)
Fruits consumption: occasionally	0.006517 (0.00271)	0.005170 (0.00271)	0.005170 (0.00271)	0.005170 (0.00271)	0.00476 (0.00271)	0.00476 (0.00271)	0.00476 (0.00271)	0.00476 (0.00271)	0.00476 (0.00271)	0.00476 (0.00271)	0.00476 (0.00271)	0.00476 (0.00271)	-0.004747 (0.00026)	-0.004747 (0.00026)	-0.004747 (0.00026)	-0.004747 (0.00026)
Fish consumption: daily	0.005151 (0.00272)	0.005114 (0.00272)	0.005114 (0.00272)	0.005114 (0.00272)	0.0117*** (0.00026)	0.0117*** (0.00026)	0.0117*** (0.00026)	0.0117*** (0.00026)	0.0117*** (0.00026)	0.0117*** (0.00026)	0.0117*** (0.00026)	0.0117*** (0.00026)	0.0117*** (0.00026)	0.0117*** (0.00026)	0.0117*** (0.00026)	0.0117*** (0.00026)
Fish consumption: weekly	0.000266 (0.00256)	0.000266 (0.00256)	0.000266 (0.00256)	0.000266 (0.00256)	0.000265 (0.00256)	0.000265 (0.00256)	0.000265 (0.00256)	0.000265 (0.00256)	0.000265 (0.00256)	0.000265 (0.00256)	0.000265 (0.00256)	0.000265 (0.00256)	-0.00367 (0.00026)	-0.00367 (0.00026)	-0.00367 (0.00026)	-0.00367 (0.00026)
Fish consumption: occasionally	0.001807 (0.00241)	0.001807 (0.00241)	0.001807 (0.00241)	0.001807 (0.00241)	0.00176 (0.00241)	0.00176 (0.00241)	0.00176 (0.00241)	0.00176 (0.00241)	0.00176 (0.00241)	0.00176 (0.00241)	0.00176 (0.00241)	0.00176 (0.00241)	-0.00481 (0.00026)	-0.00481 (0.00026)	-0.00481 (0.00026)	-0.00481 (0.00026)
Chicken/Meat consumption: daily	0.0158*** (0.00313)	0.0157*** (0.00313)	0.0157*** (0.00313)	0.0157*** (0.00313)	0.0116*** (0.00320)	0.0116*** (0.00320)	0.0116*** (0.00320)	0.0116*** (0.00320)	0.0116*** (0.00320)	0.0116*** (0.00320)	0.0116*** (0.00320)	0.0116*** (0.00320)	-0.0388*** (0.00026)	-0.0388*** (0.00026)	-0.0388*** (0.00026)	-0.0388*** (0.00026)
Chicken/Meat consumption: weekly	0.004155* (0.00277)	0.004151* (0.00277)	0.004151* (0.00277)	0.004151* (0.00277)	0.004131* (0.00300)	0.004131* (0.00300)	0.004131* (0.00300)	0.004131* (0.00300)	0.004131* (0.00300)	0.004131* (0.00300)	0.004131* (0.00300)	0.004131* (0.00300)	-0.0109*** (0.00026)	-0.0109*** (0.00026)	-0.0109*** (0.00026)	-0.0109*** (0.00026)
Chicken/Meat consumption: occasionally	0.003630 (0.00288)	-0.003137 (0.00288)	-0.003137 (0.00288)	-0.003137 (0.00288)	-0.003577 (0.00026)	-0.003577 (0.00026)	-0.003577 (0.00026)	-0.003577 (0.00026)	-0.003577 (0.00026)	-0.003577 (0.00026)	-0.003577 (0.00026)	-0.003577 (0.00026)	-0.0113*** (0.00026)	-0.0113*** (0.00026)	-0.0113*** (0.00026)	-0.0113*** (0.00026)
Fried food consumption: daily	0.002431 (0.00241)	0.002431 (0.00241)	0.002431 (0.00241)	0.002431 (0.00241)	0.002431 (0.00241)	0.002431 (0.00241)	0.002431 (0.00241)	0.002431 (0.00241)	0.002431 (0.00241)	0.002431 (0.00241)	0.002431 (0.00241)	0.002431 (0.00241)	-0.0117*** (0.00026)	-0.0117*** (0.00026)	-0.0117*** (0.00026)	-0.0117*** (0.00026)
Fried food consumption: weekly	0.002277 (0.00241)	-0.002277 (0.00241)	-0.002277 (0.00241)	-0.002277 (0.00241)	0.002431 (0.00026)	0.002431 (0.00026)	0.002431 (0.00026)	0.002431 (0.00026)	0.002431 (0.00026)	0.002431 (0.00026)	0.002431 (0.00026)	0.002431 (0.00026)	-0.0117*** (0.00026)	-0.0117*** (0.00026)	-0.0117*** (0.00026)	-0.0117*** (0.00026)
Fried food consumption: occasionally	0.002431 (0.00241)	0.002431 (0.00241)	0.002431 (0.00241)	0.002431 (0.00241)	0.002431 (0.00241)	0.002431 (0.00241)	0.002431 (0.00241)	0.002431 (0.00241)	0.002431 (0.00241)	0.002431 (0.00241)	0.002431 (0.00241)	0.002431 (0.00241)	-0.0117*** (0.00026)	-0.0117*** (0.00026)	-0.0117*** (0.00026)	-0.0117*** (0.00026)
Aerated drinks consumption: daily	0.006146* (0.00347)	0.006146* (0.00347)	0.006146* (0.00347)	0.006146* (0.00347)	0.006146* (0.00347)	0.006146* (0.00347)	0.006146* (0.00347)	0.006146* (0.00347)	0.006146* (0.00347)	0.006146* (0.00347)	0.006146* (0.00347)	0.006146* (0.00347)	0.0141*** (0.00036)	0.0141*** (0.00036)	0.0141*** (0.00036)	0.0141*** (0.00036)
Aerated drinks consumption: weekly	-0.00218 (0.00218)	-0.00218 (0.00218)	-0.00218 (0.00218)	-0.00218 (0.00218)	-0.00233 (0.00219)	-0.00233 (0.00219)	-0.00233 (0.00219)	-0.00233 (0.00219)	-0.002358 (0.00026)	-0.002358 (0.00026)	-0.002358 (0.00026)	-0.002358 (0.00026)	-0.002358 (0.00036)	-0.002358 (0.00036)	-0.002358 (0.00036)	-0.002358 (0.00036)
Aerated drinks consumption: occasionally	0.000311 (0.00158)	6.43e-05 (0.00158)	6.43e-05 (0.00158)	6.43e-05 (0.00158)	0.000308 (0.00158)	0.000308 (0.00158)	0.000308 (0.00158)	0.000308 (0.00158)	-0.000992 (0.00026)	-0.000992 (0.00026)	-0.0009					

Table A3: Effect of obesity on diastolic reading, systolic reading, blood glucose level, and haemoglobin level, under OLS, 2SLS, SLF, and SLF+IV models (without controlling for Place of residence: urban)

VARIABLES	log of average diastolic reading				log of average systolic reading				log of blood glucose level				log of haemoglobin level			
	OLS	2SLS	SLF	SLF+IV	OLS	2SLS	SLF	SLF+IV	OLS	2SLS	SLF	SLF+IV	OLS	2SLS	SLF	SLF+IV
Treatment variable: obese	0.0560*** (0.00193)	0.112*** (0.00389)	0.0560*** (0.00193)	0.0823*** (0.00199)	0.0560*** (0.00199)	0.0891*** (0.00196)	0.0560*** (0.00199)	0.0482*** (0.00199)	0.0558*** (0.00355)	0.0972*** (0.00355)	0.0588*** (0.00355)	0.0607*** (0.00355)	0.122*** (0.00241)	0.0232*** (0.00241)	0.122*** (0.00241)	0.0232*** (0.00241)
Current Age	0.0317*** (0.000201)	0.0298*** (0.000203)	0.0319*** (0.000202)	0.03039*** (0.000202)	0.03039*** (0.000201)	0.03036*** (0.000201)	0.03036*** (0.000201)	0.03036*** (0.000201)	0.02636*** (0.000201)	0.0273*** (0.000201)	0.02625*** (0.000201)	0.0273*** (0.000201)	0.02245 (0.000274)	0.02245 (0.000274)	0.02245 (0.000274)	0.02245 (0.000274)
Alcohol consumption: yes	0.0276*** (0.000128)	0.0284*** (0.000128)	0.0284*** (0.000127)	0.0289*** (0.000127)	0.0289*** (0.000127)	0.0376*** (0.000127)	0.0376*** (0.000127)	0.0377*** (0.000127)	0.0378*** (0.000127)	0.0378*** (0.000127)	0.0378*** (0.000127)	0.0378*** (0.000127)	0.0225 (0.000274)	0.0225 (0.000274)	0.0225 (0.000274)	0.0225 (0.000274)
Milk/Curd consumption: daily	-0.009385 (0.00210)	-0.009188 (0.00210)	-0.00917 (0.00210)	-0.009175 (0.00210)	-0.009175 (0.00210)	-0.009175 (0.00210)	-0.009175 (0.00210)	-0.009175 (0.00210)	-0.009175 (0.00210)	-0.009175 (0.00210)	-0.009175 (0.00210)	-0.009175 (0.00210)	0.00530 (0.00514)	0.00530 (0.00514)	0.00530 (0.00514)	0.00530 (0.00514)
Milk/Curd consumption: weekly	0.00666 (0.00219)	0.006541 (0.00219)	0.006541 (0.00219)	0.006541 (0.00219)	0.006541 (0.00219)	0.007540 (0.00219)	0.007540 (0.00219)	0.007540 (0.00219)	0.007540 (0.00219)	0.007540 (0.00219)	0.007540 (0.00219)	0.007540 (0.00219)	0.00514 (0.00514)	0.00514 (0.00514)	0.00514 (0.00514)	0.00514 (0.00514)
Milk/Curd consumption: occasionally	0.00168 (0.00216)	0.000635 (0.00216)	0.000635 (0.00216)	0.000635 (0.00216)	0.000635 (0.00216)	2.15e-05 (0.00216)	0.00521*** (0.00216)	0.00521*** (0.00216)	0.00586*** (0.00216)	0.00586*** (0.00216)	0.00586*** (0.00216)	0.00586*** (0.00216)	-0.00949 (0.00514)	-0.00949 (0.00514)	-0.00949 (0.00514)	-0.00949 (0.00514)
Pulses/Beans consumption: daily	0.0255*** (0.0125)	0.0254*** (0.0125)	0.0254*** (0.0125)	0.0254*** (0.0125)	0.0254*** (0.0125)	0.0254*** (0.0125)	0.0254*** (0.0125)	0.0254*** (0.0125)	0.0254*** (0.0125)	0.0254*** (0.0125)	0.0254*** (0.0125)	0.0254*** (0.0125)	0.00556 (0.00514)	0.00556 (0.00514)	0.00556 (0.00514)	0.00556 (0.00514)
Pulses/Beans consumption: weekly	0.0296*** (0.0125)	0.0289*** (0.0125)	0.0289*** (0.0125)	0.0289*** (0.0125)	0.0289*** (0.0125)	0.0289*** (0.0125)	0.0289*** (0.0125)	0.0289*** (0.0125)	0.0289*** (0.0125)	0.0289*** (0.0125)	0.0289*** (0.0125)	0.0289*** (0.0125)	0.00565 (0.00514)	0.00565 (0.00514)	0.00565 (0.00514)	0.00565 (0.00514)
Pulses/Beans consumption: occasionally	0.0321*** (0.0126)	0.0321*** (0.0126)	0.0321*** (0.0126)	0.0321*** (0.0126)	0.0321*** (0.0126)	0.0321*** (0.0126)	0.0321*** (0.0126)	0.0321*** (0.0126)	0.0321*** (0.0126)	0.0321*** (0.0126)	0.0321*** (0.0126)	0.0321*** (0.0126)	0.0130 (0.00514)	0.0130 (0.00514)	0.0130 (0.00514)	0.0130 (0.00514)
Dark Green Leafy Vegetables consumption: daily	0.06550 (0.0204)	0.06341 (0.0204)	0.06341 (0.0204)	0.06341 (0.0204)	0.06341 (0.0204)	0.00435 (0.0169)	0.00435 (0.0169)	0.00435 (0.0169)	0.00435 (0.0169)	0.00435 (0.0169)	0.00435 (0.0169)	0.00435 (0.0169)	-0.0131 (0.0169)	-0.0131 (0.0169)	-0.0131 (0.0169)	-0.0131 (0.0169)
Dark Green Leafy Vegetables consumption: weekly	0.00881 (0.0204)	-0.00755 (0.0204)	-0.00755 (0.0204)	-0.00755 (0.0204)	-0.00755 (0.0204)	-1.16e-05 (0.0204)	-0.0178 (0.0204)	-0.0178 (0.0204)	-0.0178 (0.0204)	-0.0178 (0.0204)	-0.0178 (0.0204)	-0.0178 (0.0204)	-0.0124 (0.0169)	-0.0124 (0.0169)	-0.0124 (0.0169)	-0.0124 (0.0169)
Dark Green Leafy Vegetables consumption: occasionally	-0.00519 (0.0204)	-0.00519 (0.0204)	-0.00519 (0.0204)	-0.00519 (0.0204)	-0.00519 (0.0204)	-0.00519 (0.0204)	-0.00519 (0.0204)	-0.00519 (0.0204)	-0.00519 (0.0204)	-0.00519 (0.0204)	-0.00519 (0.0204)	-0.00519 (0.0204)	0.00617 (0.0169)	0.00617 (0.0169)	0.00617 (0.0169)	0.00617 (0.0169)
Fruits consumption: daily	-0.00811 (0.00534)	-0.00787 (0.00534)	-0.00787 (0.00534)	-0.00787 (0.00534)	-0.00787 (0.00534)	-0.0127*** (0.00534)	-0.0128*** (0.00534)	-0.0128*** (0.00534)	-0.0128*** (0.00534)	-0.0128*** (0.00534)	-0.0128*** (0.00534)	-0.0128*** (0.00534)	-0.00627 (0.00534)	-0.00627 (0.00534)	-0.00627 (0.00534)	-0.00627 (0.00534)
Fruits consumption: weekly	0.00238 (0.0120)	0.00238 (0.0120)	0.00238 (0.0120)	0.00238 (0.0120)	0.00238 (0.0120)	0.00238 (0.0120)	0.00238 (0.0120)	0.00238 (0.0120)	0.00238 (0.0120)	0.00238 (0.0120)	0.00238 (0.0120)	0.00238 (0.0120)	-0.0139 (0.0169)	-0.0139 (0.0169)	-0.0139 (0.0169)	-0.0139 (0.0169)
Fruits consumption: occasionally	0.00369 (0.0125)	0.00375 (0.0125)	0.00375 (0.0125)	0.00375 (0.0125)	0.00375 (0.0125)	0.00375 (0.0125)	0.00375 (0.0125)	0.00375 (0.0125)	0.00375 (0.0125)	0.00375 (0.0125)	0.00375 (0.0125)	0.00375 (0.0125)	0.00366 (0.00514)	0.00366 (0.00514)	0.00366 (0.00514)	0.00366 (0.00514)
Fish consumption: daily	-0.0119*** (0.0124)	-0.0121*** (0.0124)	-0.0121*** (0.0124)	-0.0121*** (0.0124)	-0.0121*** (0.0124)	-0.00851*** (0.0124)	-0.00851*** (0.0124)	-0.00851*** (0.0124)	-0.00851*** (0.0124)	-0.00851*** (0.0124)	-0.00851*** (0.0124)	-0.00851*** (0.0124)	-0.0124 (0.0169)	-0.0124 (0.0169)	-0.0124 (0.0169)	-0.0124 (0.0169)
Fish consumption: weekly	0.00434 (0.0253)	0.00434 (0.0253)	0.00434 (0.0253)	0.00434 (0.0253)	0.00434 (0.0253)	0.00434 (0.0253)	0.00434 (0.0253)	0.00434 (0.0253)	0.00434 (0.0253)	0.00434 (0.0253)	0.00434 (0.0253)	0.00434 (0.0253)	-0.0118 (0.0169)	-0.0118 (0.0169)	-0.0118 (0.0169)	-0.0118 (0.0169)
Fish consumption: occasionally	0.005050 (0.0125)	0.005156 (0.0125)	0.005156 (0.0125)	0.005156 (0.0125)	0.005156 (0.0125)	0.008865 (0.0125)	0.008865 (0.0125)	0.008865 (0.0125)	0.008865 (0.0125)	0.008865 (0.0125)	0.008865 (0.0125)	0.008865 (0.0125)	-0.15e-05 (0.0169)	-0.15e-05 (0.0169)	-0.15e-05 (0.0169)	-0.15e-05 (0.0169)
Chicken/Meat consumption: daily	0.0154*** (0.0125)	0.0148*** (0.0125)	0.0148*** (0.0125)	0.0148*** (0.0125)	0.0148*** (0.0125)	-0.00121*** (0.0125)	-0.00121*** (0.0125)	-0.00121*** (0.0125)	-0.00121*** (0.0125)	-0.00121*** (0.0125)	-0.00121*** (0.0125)	-0.00121*** (0.0125)	-0.00121*** (0.0125)	-0.00121*** (0.0125)	-0.00121*** (0.0125)	-0.00121*** (0.0125)
Chicken/Meat consumption: weekly	0.001419 (0.00256)	0.001406 (0.00256)	0.001406 (0.00256)	0.001406 (0.00256)	0.001406 (0.00256)	0.003666 (0.00256)	0.003666 (0.00256)	0.003666 (0.00256)	0.003666 (0.00256)	0.003666 (0.00256)	0.003666 (0.00256)	0.003666 (0.00256)	-0.0118*** (0.0169)	-0.0118*** (0.0169)	-0.0118*** (0.0169)	-0.0118*** (0.0169)
Chicken/Meat consumption: occasionally	0.001449*** (0.00246)	0.004585*** (0.00246)	0.004585*** (0.00246)	0.004585*** (0.00246)	0.004585*** (0.00246)	0.003522 (0.00246)	0.003522 (0.00246)	0.003522 (0.00246)	0.003522 (0.00246)	0.003522 (0.00246)	0.003522 (0.00246)	0.003522 (0.00246)	-0.0175*** (0.0169)	-0.0175*** (0.0169)	-0.0175*** (0.0169)	-0.0175*** (0.0169)
Aerated food consumption: daily	0.003340 (0.00343)	0.003432 (0.00343)	0.003432 (0.00343)	0.003432 (0.00343)	0.003432 (0.00343)	0.0149*** (0.00343)	0.0149*** (0.00343)	0.0149*** (0.00343)	0.0149*** (0.00343)	0.0149*** (0.00343)	0.0149*** (0.00343)	0.0149*** (0.00343)	-0.00403 (0.0169)	-0.00403 (0.0169)	-0.00403 (0.0169)	-0.00403 (0.0169)
Aerated food consumption: weekly	-0.002322 (0.00229)	-0.002322 (0.00229)	-0.002322 (0.00229)	-0.002322 (0.00229)	-0.002322 (0.00229)	-0.002322 (0.00229)	-0.002322 (0.00229)	-0.002322 (0.00229)	-0.002322 (0.00229)	-0.002322 (0.00229)	-0.002322 (0.00229)	-0.002322 (0.00229)	-0.00118*** (0.0169)	-0.00118*** (0.0169)	-0.00118*** (0.0169)	-0.00118*** (0.0169)
Aerated food consumption: occasionally	-0.002307 (0.00229)	-0.002307 (0.00229)	-0.002307 (0.00229)	-0.002307 (0.00229)	-0.002307 (0.00229)	-0.00330 (0.00229)	-0.00330 (0.00229)	-0.00330 (0.00229)	-0.00330 (0.00229)	-0.00330 (0.00229)	-0.00330 (0.00229)	-0.00330 (0.00229)	-0.00122*** (0.0169)	-0.00122*** (0.0169)	-0.00122*** (0.0169)	-0.00122*** (0.0169)
Aerated drinks consumption: daily	0.00289 (0.02417)	0.00289 (0.02417)	0.00289 (0.02417)	0.00289 (0.02417)	0.00289 (0.02417)	-0.00225 (0.02417)	-0.00225 (0.02417)	-0.00225 (0.02417)	-0.00225 (0.02417)	-0.00225 (0.02417)	-0.00225 (0.02417)	-0.00225 (0.02417)	-0.00118*** (0.0169)	-0.00118*** (0.0169)	-0.00118*** (0.0169)	-0.00118*** (0.0169)
Aerated drinks consumption: weekly	4.215*** (0.0243)	4.215*** (0.0243)	4.215*** (0.0243)	4.215*** (0.0243)	4.215*** (0.0243)	4.623*** (0.0243)	4.623*** (0.0243)	4.623*** (0.0243)	4.623*** (0.0243)	4.623*** (0.0243)	4.623*** (0.0243)	4.623*** (0.0243)	4.571*** (0.0243)	4.571*** (0.0243)	4.571*** (0.0243)	4.571*** (0.0243)
Aerated drinks consumption: occasionally	4.215*** (0.0243)	4.215*** (0.0243)	4.215*** (0.0243)	4.215*** (0.0243)	4.215*** (0.0243)	4.623*** (0.0243)	4.623*** (0.0243)	4.623*** (0.0243)	4.623*** (0.0243)	4.623*** (0.0243)	4.623*** (0.0243)	4.623*** (0.0243)	4.568*** (0.0243)	4.568*** (0.0243)	4.568*** (0.0243)	4.568*** (0.0243)
Latent factor																
Constant	39.105	39.101	42.120	42.074	39.106	39.102	42.120	42.074	41.552	41.475	42.120	42.074	41.411	41.411	42.120	42.074

Notes:

Robust standard errors in the parentheses

** $p < 0.01$, * $p < 0.05$, $^{+}p < 0.1$