# *Erratum*

# Authors' corrigenda/corrections des auteurs on synthetic data method to incorporate external information into a current study

Tian GU[1]*, Jeremy M. G. TAYLOR[1], Wenting CHENG[1], and Bhramar MUKHERJEE[1]

[1]*Department of Biostatistics, University of Michigan, Ann Arbor, MI 48105, U.S.A.*

In Gu, Taylor, Cheng & Mukherjee (2019), four notations used in Appendices A1, A2 and A3, namely $\Omega, \omega^2, \sigma^2$ and $\eta^2$, should have been replaced by some of the existing notations. The correct notations that they correspond to are listed as follows:

- $\Omega = \sigma_X^2$, the variance of X;
- $\omega^2 = \sigma_\beta^2$, as introduced in equation (5) in the main paper;
- $\sigma^2 = \sigma_\gamma^2$, as introduced in equation (6) in the main paper;
- $\eta^2 = \sigma_\theta^2$, as introduced in equation (7) in the main paper.

Note that all the results in Appendices A1, A2 and A3 were correct and were not affected by the notation correction. The corrected Appendices A1, A2 and A3 are attached below.

## REFERENCE

Gu, T., Taylor, J. M. G., Cheng, W., & Mukherjee, B. (2019). Synthetic data method to incorporate external information into a current study. *The Canadian Journal of Statistics*, 47, 580–603.

## APPENDIX

*Derivation of asymptotic variances for the special case 1*

## Appendix A1. Approach 1: Synthetic Data Method

If the synthetic data approach is applied, and under the assumption that the true value of $\beta$ and $\sigma_\beta$ are used to generate the synthetic data, then the combined data will have the same distribution as a dataset of size n+m in which m values of B have been removed. For this particular data structure, it is possible to obtain formulas for the asymptotic variance of the maximum likelihood estimate (MLE) of $\gamma$. In particular, Gourieroux & Monfort (1981) gave the exact expression of the MLE and the corresponding asymptotic covariance in such case. The likelihood for the combined data is $\prod_{i=1}^{n} f(Y_i, B_i|X_i) \times \prod_{i=n+1}^{n+m} f(Y_i|X_i)$, which can be

---

* *Corresponding author: gutian@umich.edu*

rewritten as $\prod_{i=1}^{n+m} f(Y_i|X_i) \times \prod_{i=1}^{n} f(B_i|X_i, Y_i)$. Based on this likelihood, they introduced a set of transformed parameters, and re-parameterized the distributions (5)–(7). They then identified the 1-to-1 relationship among the original parameters and the new set of parameters, which we will explain in the subsequent paragraph.

We obtain the estimators of the original parameters by the re-parameterization method, and then apply the delta method to get the asymptotic variance of $\hat{\gamma}_B$ and $\hat{\gamma}_X$. According to Gourieroux & Monfort (1981), we introduce a set of transformed parameters $a, b, c, d$ and $e$, and re-parameterize the distributions (5)–(7) as $Y|X \sim N(bX, a^2)$, and $B|Y, X \sim N(dY + eX, c^2)$. Then we identify the 1-to-1 relationship among the original parameters and the new set of parameters:

$$a^2 = \sigma_\gamma^2 + \gamma_B^2 \sigma_\theta^2$$

$$b = \gamma_X + \theta \gamma_B$$

$$c^2 = \frac{\sigma_\gamma^2 \sigma_\theta^2}{a^2}$$

$$d = \frac{\gamma_B \sigma_\theta^2}{a^2}$$

$$e = \theta - db. \tag{1}$$

The MLE of $a$, $b$ and their asymptotic variances are easy to obtain from the linear model $Y_i = bX_i + u_i, Var(u_i) = a^2$, where i = 1,…, n+m. Similarly, the MLE of $c$, $d$ and $e$ and their asymptotic variances are easy to obtain from the linear model $B_i = dY_i + eX_i + v_i, Var(v_i) = c^2$ where i = 1,…,n. The estimators of the original parameters are obtained through the relationship derived from Equation (1), where

$$\theta = bd + e$$

$$\sigma_\theta^2 = a^2 d^2 + c^2$$

$$\gamma_B = \frac{a^2 d}{\sigma_\theta^2}$$

$$\gamma_X = b - \gamma_B \theta$$

$$\sigma_\gamma^2 = \frac{a^2 c^2}{\sigma_\theta^2}$$

and the asymptotic variance of $\hat{\gamma}_X$ and $\hat{\gamma}_B$ can be derived using the delta method:

$$\begin{cases} Var(\hat{\gamma}_B) = \frac{1}{n}\left[ \frac{\sigma_\gamma^2}{\sigma_\theta^2} + 2(\lambda - 1)\frac{\gamma_B^2 \sigma_\gamma^4}{\sigma_\beta^4} \right] \\ Var(\hat{\gamma}_X) = \theta^2 Var(\hat{\gamma}_B) + \frac{1}{n}\frac{\sigma_\gamma^2}{\sigma_X^2}\frac{\lambda\sigma_\gamma^2 + \gamma_B^2 \sigma_\theta^2}{\sigma_\beta^2}. \end{cases}$$

Therefore, we find the relative efficiency gain of $Var(\hat{\gamma}) = Var(\hat{\gamma}_X, \hat{\gamma}_B)^T$ by adding m synthetic data observations compared to the original dataset of size n is

$$ARE[Var(\hat{\gamma})] = 1 - (1 - \lambda)\left( \frac{\frac{2\sigma_\gamma^2 \gamma_B^2 \sigma_\theta^2}{\sigma_\beta^4} + \frac{\sigma_\theta^2}{\sigma_\theta^2 + \sigma_X^2 \theta^2}\frac{\sigma_\gamma^2(2\sigma_\gamma^2 - \sigma_\beta^2)}{\sigma_\beta^4}}{\frac{2\gamma_B^2 \sigma_\theta^2 \sigma_\gamma^2}{\sigma_\beta^4}} \right),$$

where $\theta = \frac{\beta - \gamma_X}{\gamma_B}$ and $\sigma_\theta^2 = \frac{\sigma_\beta^2 - \sigma_\gamma^2}{\gamma_B^2}$. When m gets very large such that $\lambda \approx 0$, ARE[Var($\hat{\gamma}_X$)] =
$1 - \frac{\sigma_\theta^2}{\sigma_\theta^2 + \sigma_X^2 \theta^2} \frac{\sigma_\gamma^2}{\sigma_\gamma^2 + \gamma_B^2 \sigma_\theta^2} \frac{2\sigma_\gamma^2 - \sigma_\beta^2}{\sigma_\beta^2} - \frac{2\sigma_\gamma^2 \gamma_B^2 \sigma_\theta^2}{\sigma_\beta^4}$, and ARE[Var($\hat{\gamma}_B$)] $= 1 - \frac{2\sigma_\gamma^2 \gamma_B^2 \sigma_\theta^2}{\sigma_\beta^4}$. This demonstrates some gain in efficiency for both $\gamma_X$ and $\gamma_B$.

## Appendix A2. Approach 2: Constrained MLE

Depending on the information available from the external model Y|X, two possible situations correspond to two different constraints:

- *Approach 2.1: Only the estimated coefficient $\beta$ is known from model (5)*
  From model (5)–(7), it is easy to see that the constraint is $\theta = \frac{\beta - \gamma_X}{\gamma_B}$, describing the relationship between the unknown variable $\theta$, the known variable $\beta$ and the target variable $\gamma$. The log-likelihood is given by

$$l = l(\gamma, \theta, \sigma_\gamma^2, \sigma_\theta^2)$$

$$= \sum_{i=1}^{n} logf(Y_i, B_i | X_i) = \sum_{i=1}^{n} logf(Y_i | X_i, B_i; \gamma, \sigma_\gamma^2) + \sum_{i=1}^{n} logf(B_i | X_i; \theta, \sigma_\theta^2)$$

$$= -\frac{n}{2}log(\sigma_\gamma^2) - \frac{1}{2\sigma_\gamma^2} \sum_{i=1}^{n} \left(Y_i - \gamma_X X_i - \gamma_B B_i\right)^2 - \frac{n}{2}log(\sigma_\theta^2) - \frac{1}{2\sigma_\theta^2} \sum_{i=1}^{n} \left(B_i - \theta X_i\right)^2. \quad (2)$$

The goal is to maximize the log-likelihood (2) over $\gamma$, $\sigma_\gamma$ and $\sigma_\theta$ subject to the constraint $\theta = \theta^* = \frac{\beta^* - \gamma_X}{\gamma_B}$. By replacing $\theta$ with $\theta^*$, taking the second derivative over $\gamma$, and taking the inverse of the information matrix, we obtain the asymptotic variance of $\hat{\gamma}$:

$$\text{Var}(\hat{\gamma}) = \mathbf{I}^{-1} = \frac{1}{n} \frac{\sigma_\gamma^2}{\sigma_\theta^2} \begin{pmatrix} \theta^{*2} + \frac{\sigma_\theta^4 \gamma_B^2}{\sigma_\gamma^2 + \sigma_\theta^2 \gamma_B^2} \sigma_X^{-2} & -\theta^* \\ -\theta^* & 1 \end{pmatrix}. \quad (3)$$

Thus, the ARE of Var($\hat{\gamma}$) from the constrained MLE compared to the standard MLE is

$$\text{ARE}[\text{Var}(\hat{\gamma})] = \begin{pmatrix} 1 - \frac{\sigma_\theta^2}{\sigma_\theta^2 + \sigma_X^2 \theta^{*2}} \frac{\sigma_\gamma^2}{\sigma_\gamma^2 + \gamma_B^2 \sigma_\theta^2} \\ 1 \end{pmatrix},$$

where we notice that there is some gain in efficiency for $\gamma_X$ but no gain in efficiency for $\gamma_B$. We can see that the largest gain in efficiency is when $\gamma_B$, $\theta$ and $\sigma_X$ are small.

- *Approach 2.2: Both of the estimated coefficient $\beta$ and the standard deviation $\sigma_\beta$ are known from model (5)*
  In this situation, knowing the true $\sigma_\beta$ gives us more information which is incorporated through an additional constraint. In addition to the constraint $\theta = \theta^* = \frac{\beta - \gamma_X}{\gamma_B}$ derived in approach 2.1, we add another constraint Var(Y|X) $= \sigma_\beta^2 = \gamma_B^2 \sigma_\theta^2 + \sigma_\gamma^2$, where $\sigma_\theta^2 = \sigma_\theta^{*2} = \frac{\sigma_\beta^{*2} - \sigma_\gamma^2}{\gamma_B^2}$. Then we maximize log-likelihood (2) with respect to $\gamma$ and $\sigma_\gamma^2$ at fixed values $\sigma_\theta^2 = \sigma_\theta^{*2}, \theta = \theta^*$. Note that different from approach 2.1, $\sigma_\gamma^2$ and $\gamma$ are not independent anymore. Thus, we need to

consider $\sigma_\gamma^2$ in the information matrix, and take the inverse of a 3×3 matrix to get the correct asymptotic variance. Let $\boldsymbol{\phi} = (\gamma, \sigma_\gamma^2)^{\mathrm{T}}$,

$$
\mathbf{I} = -\mathrm{E}_{\mathrm{XB}}\left(\frac{\partial^2 l}{\partial \boldsymbol{\phi}\partial\boldsymbol{\phi}^{\mathrm{T}}}\right)
$$

$$
= n\begin{pmatrix}
\left(\frac{1}{\sigma_\gamma^2} + \frac{1}{\gamma_{\mathrm{B}}^2\sigma_\theta^{*2}}\right)\sigma_{\mathrm{X}}^2 & \left(\frac{1}{\sigma_\gamma^2} + \frac{1}{\gamma_{\mathrm{B}}^2\sigma_\theta^{*2}}\right)\sigma_{\mathrm{X}}^2\theta^* & 0 \\
\left(\frac{1}{\sigma_\gamma^2} + \frac{1}{\gamma_{\mathrm{B}}^2\sigma_\theta^{*2}}\right)\sigma_{\mathrm{X}}^2\theta^* & \left(\frac{1}{\sigma_\gamma^2} + \frac{1}{\gamma_{\mathrm{B}}^2\sigma_\theta^{*2}}\right)\left(\sigma_\theta^{*2} + \sigma_{\mathrm{X}}^2\theta^{*2}\right) + \frac{1}{\gamma_{\mathrm{B}}^2} & \frac{1}{\sigma_\theta^{*2}\gamma_{\mathrm{B}}^3} \\
0 & \frac{1}{\sigma_\theta^{*2}\gamma_{\mathrm{B}}^3} & \frac{1}{2}\left(\frac{1}{\sigma_\gamma^4} + \frac{1}{\gamma_{\mathrm{B}}^4\sigma_\theta^{*4}}\right)
\end{pmatrix}.
$$

By taking the inverse of $\mathbf{I}$, we can get the asymptotic variance of $\hat{\boldsymbol{\gamma}}$:

$$
\begin{cases}
\mathrm{Var}(\hat{\gamma}_{\mathrm{B}}) = \frac{1}{n}\frac{\sigma_\gamma^2}{\sigma_\theta^{*2}}\frac{\sigma_\gamma^4 + \gamma_{\mathrm{B}}^4\sigma_\theta^{*4}}{\left(\sigma_\gamma^2 + \gamma_{\mathrm{B}}^2\sigma_\theta^{*2}\right)^2} = \frac{1}{n}\frac{\gamma_{\mathrm{B}}^2\sigma_\gamma^2}{\sigma_\beta^{*2} - \sigma_\gamma^2}\frac{\sigma_\gamma^4 + \left(\sigma_\beta^{*2} - \sigma_\gamma^2\right)^2}{\sigma_\beta^{*4}} \\[2ex]
\mathrm{Var}(\hat{\gamma}_{\mathrm{X}}) = \frac{1}{n}\frac{\sigma_\gamma^2}{\sigma_\theta^{*2}}\frac{1}{\left(\sigma_\gamma^2 + \gamma_{\mathrm{B}}^2\sigma_\theta^{*2}\right)^2}\left[\left(\sigma_\theta^{*2}\sigma_{\mathrm{X}}^{-2} + \theta^{*2}\right)\left(\sigma_\gamma^4 + \gamma_{\mathrm{B}}^4\sigma_\theta^{*4}\right) - \left(\sigma_\gamma^2 - \gamma_{\mathrm{B}}^2\sigma_\theta^{*2}\right)\sigma_\gamma^2\sigma_\theta^{*2}\sigma_{\mathrm{X}}^{-2}\right] \\[2ex]
\quad = \left(\sigma_\theta^{*2}\sigma_{\mathrm{X}}^{-2} + \theta^{*2}\right)\mathrm{Var}(\hat{\gamma}_{\mathrm{B}}) - \frac{1}{n}\sigma_\gamma^4\sigma_{\mathrm{X}}^{-2}\frac{\sigma_\gamma^2 - \gamma_{\mathrm{B}}^2\sigma_\theta^{*2}}{\sigma_\beta^{*4}}
\end{cases}.
$$

Thus, we can obtain the identical ARE to the synthetic data method (approach 1). This demonstrates the asymptotic equivalence of the synthetic data approach with large m compared to the constrained ML approach that uses knowledge of all the parameters in the Y|X distribution.

## Appendix A3. Approach 3: Constrained Semiparametric MLE

This approach assumes that $\beta$ is known, but does not assume that $\sigma_\beta$ is known. To calculate the asymptotic variance of $\hat{\boldsymbol{\gamma}}$ in this approach, we need three matrices $\mathbf{I}$, $\mathbf{C}$ and $\mathbf{L}$. After some algebra, it can be shown that $\mathbf{C} = \frac{\sigma_{\mathrm{X}}^2}{\sigma_\beta^2}(1, \theta^*)^{\mathrm{T}}, \mathbf{L} = \frac{n\gamma_{\mathrm{B}}^2\sigma_\theta^*\sigma_{\mathrm{X}}^2}{\sigma_\beta^{*4}}$. Thus,

$$
\mathrm{Cov}(\hat{\boldsymbol{\gamma}}) = (\mathbf{I} + \mathbf{C}\mathbf{L}^{-1}\mathbf{C}^{\mathrm{T}})^{-1} = \frac{1}{n}\frac{\sigma_\gamma^2}{\sigma_\theta^2}\begin{pmatrix}
\theta^{*2} + \frac{\sigma_\theta^4\gamma_{\mathrm{B}}^2}{\sigma_\gamma^2 + \sigma_\theta^2\gamma_{\mathrm{B}}^2}\sigma_{\mathrm{X}}^{-2} & -\theta^* \\
-\theta^* & 1
\end{pmatrix},
$$

which is identical to approach 2.1.