This is a work of the United States Government. In accordance with 17 U.S.C. 105, no copyright protection is available for such works under U.S. Law.

Public Domain Mark 1.0 https://creativecommons.org/publicdomain/mark/1.0/

Access to this work was provided by the University of Maryland, Baltimore County (UMBC) ScholarWorks@UMBC digital repository on the Maryland Shared Open Access (MD-SOAR) platform.

## Please provide feedback

Please support the ScholarWorks@UMBC repository by emailing <u>scholarworks-</u> <u>group@umbc.edu</u> and telling us what having access to this work means to you and why it's important to you. Thank you.

#### RESEARCH REPORT SERIES (Statistics #2022-02)

## Bayesian Analysis of Multiply Imputed Synthetic Data Under the Multiple Linear Regression Model

Abhishek Guin<sup>1</sup>, Anindya Roy<sup>1,2</sup>, Bimal Sinha<sup>1,2</sup>

<sup>1</sup>Department of Mathematics and Statistics, University of Maryland Baltimore County; <sup>2</sup>Center for Statistical Research and Methodology, U.S. Census Bureau

> Center for Statistical Research & Methodology Research and Methodology Directorate U.S. Census Bureau Washington, D.C. 20233

Report Issued: April 4, 2022

*Disclaimer:* This report is released to inform interested parties of research and to encourage discussion. The views expressed are those of the authors and not those of the U.S. Census Bureau.

# Bayesian Analysis of Multiply Imputed Synthetic Data under the Multiple Linear Regression Model

Abhishek Guin<sup>1</sup>, Anindya Roy<sup>1,2</sup> and Bimal Sinha<sup>1,2,\*</sup> <sup>1</sup>University of Maryland Baltimore County, <sup>2</sup>U.S. Census Bureau

#### Abstract

In this paper we consider Bayesian inference of model parameters in a multiple linear regression model when the response variable is sensitive and the covariates are not, analysis being carried out based on multiple synthetic versions of the response variable. Two scenarios of synthetic data generation are considered - plug-in sampling method and posterior predictive sampling method. We also consider the case when part of the response is sensitive and describe how to carry out full Bayesian analysis based on multiply imputed data. **Keywords: Credible sets, Partially sensitive response, Privacy** 

1 Introduction

For many statistical agencies such as the US Census Bureau it is customary to publish statistical analysis of results collected from surveys as well as the original raw microdata so others can either reproduce the analysis results or can do some further statistical analysis depending on their purpose. When the response data is sensitive or confidential, a redirect release of such data is not possible and statistical agencies often release what is known as a synthetic version of the original microdata. Fortunately, there are several ways to accomplish this goal and valid statistical analysis of such synthetic data is also quite often possible. In fact, there is a rich literature addressing methods to generate synthetic data and also the appropriate statistical methods to analyze such synthetic data, primarily based on a parametric model which is believed to generate the original data [7, 6, 5, 4, 3, 2]. We mention in passing that in this paper we deal with what is known as *partially synthetic data* rather than fully synthetic data [8].

Under the assumption of a classical linear regression model (single or multiple), valid inference about the regression coefficients and residual variance based on synthetic data has been successfully established based on both single and multiple imputations [5, 4]. Broadly, two types of synthetic data generation schemes have been considered - plug-in sampling method and posterior predictive sampling method. Inference about the regression coefficients from a frequentist point of view was extensively studied in [5, 4] and [2], based on both single imputation and multiple imputations. In

<sup>\*</sup>Center for Statistical Research and Methodology, U.S. Census Bureau, Washington, DC 20233, USA, mailto:sinha@umbc.edu

*Disclaimer*: This article is released to inform interested parties of ongoing research and to encourage discussion of work in progress. The views expressed are those of the authors and not those of the U.S. Census Bureau.

this paper we investigate the inference problem from a Bayesian point of view under the multiple linear regression model with multiple synthetic versions of the original data, thus extending the results in [1] who studied the same problem under a single imputation.

Section 2 contains a brief description of synthetic data generation under plug-in sampling and posterior predictive sampling for a very general scenario. In Section 3 we derive Bayesian inferential results for the model parameters (regression coefficients and residual variance) under the plug-in sampling method while Section 4 contains similar inferential results under the posterior predictive sampling method. In Section 5 we consider the case when only a part of the response vector is sensitive (called partially sensitive data) and provide details about how to carry out appropriate inference under both the data generation schemes.

## 2 Generating Synthetic Data

We consider two ways of generating the  $m \ge 1$  synthetic copies of the original data namely, *plug-in* sampling and posterior predictive sampling. In the former method, parameter estimates are plugged in the model to generate synthetic data. In the latter one, posterior draws of the parameter are generated using an imputed prior, which are then fed into the original model to get synthetic data.

**Plug-in Sampling**. The basic mechanism for generating synthetic data via *plug-in sampling* (PIS) is described as follows: let  $\mathbf{Y} = (\mathbf{y}_1, \ldots, \mathbf{y}_n)$  be the original confidential data, which are jointly distributed according to the probability density function (pdf)  $f_{\boldsymbol{\theta}}(\mathbf{Y})$ , where  $\boldsymbol{\theta}$  is the unknown (scalar or vector) parameter. To generate partially synthetic data, let  $\hat{\boldsymbol{\theta}} = \hat{\boldsymbol{\theta}}(\mathbf{Y})$  be the observed value of a point estimator of  $\boldsymbol{\theta}$ , and we plug it into the joint pdf of  $\mathbf{Y}$ . The resulting pdf, with the unknown  $\boldsymbol{\theta}$  replaced by the observed value  $\hat{\boldsymbol{\theta}}(\mathbf{Y})$  of the point estimator, is denoted by  $f_{\hat{\boldsymbol{\theta}}}$ . The singly imputed synthetic data, denoted by  $\mathbf{Z}$ , are then generated by drawing  $\mathbf{Z}$  from the joint pdf  $\hat{f}_{\hat{\boldsymbol{\theta}}}$ . For the multiple imputation case, we draw m > 1 samples  $\mathbf{Z}_1, \ldots, \mathbf{Z}_m$  independently from  $f_{\hat{\boldsymbol{\theta}}}$ .

**Posterior Predictive Sampling**. An alternative method to generate partially synthetic data is to use *posterior predictive sampling* (PPS) which proceeds as follows: suppose that  $\mathbf{Y} = (\mathbf{y}_1, \ldots, \mathbf{y}_n)$ are the original data which are jointly distributed according to the pdf  $f_{\theta}(\mathbf{Y})$ , where  $\theta$  is the unknown (scalar or vector) parameter. Assume a prior  $\pi(\theta)$  for  $\theta$ , then the posterior distribution of  $\theta$  given  $\mathbf{Y}$  is obtained as  $\pi(\theta \mid \mathbf{Y}) \propto \pi(\theta) f_{\theta}(\mathbf{Y})$ , and used to draw  $m \geq 1$  replications  $\theta_1^*, \ldots, \theta_m^*$ (known as posterior draws). Next, for each posterior draw of  $\theta$ , a corresponding replicate of  $\mathbf{Y}$  is generated, namely  $\mathbf{Z}_j = (\mathbf{z}_{j1}, \ldots, \mathbf{z}_{jn})'$  drawn from the pdf  $f_{\theta_i^*}(\mathbf{X})$  independently for  $j = 1, \ldots, m$ .

We conclude this chapter with an observation regarding the existence of *sufficient statistics* in the context of synthetic data. The proof is given in [2].

Lemma 2.1. Suppose that when the original data  $\boldsymbol{Y}$  are observed,  $T(\boldsymbol{Y})$  is a sufficient statistic for  $\boldsymbol{\theta}$ . Then when the synthetic data  $\boldsymbol{\mathfrak{Z}} = (\boldsymbol{Z}_1, \ldots, \boldsymbol{Z}_M)$  are observed,  $(T(\boldsymbol{Z}_1), \ldots, T(\boldsymbol{Z}_M))$  is jointly sufficient for  $\boldsymbol{\theta}$ . Furthermore, if M = 1, the sufficient statistic is simply  $T(\boldsymbol{Z}_1)$ , and if M > 1, then  $\sum_{i=1}^{M} T(\boldsymbol{Z}_i)$  is sufficient if  $f_{\boldsymbol{\theta}}(\boldsymbol{Y}) = h(\boldsymbol{Y})\psi(\boldsymbol{\theta}) \exp{\{\gamma(\boldsymbol{\theta})'T(\boldsymbol{Y})\}}$ , i.e., if  $f_{\boldsymbol{\theta}}(\boldsymbol{Y})$  belongs to the exponential family.

#### 3 Plug In Sampling method

Consider a standard multiple linear regression (MLR) model involving a sensitive response variable y and a  $p \times 1$  dimensional vector of non-sensitive predictors x. We want to generate synthetic data  $\boldsymbol{z}_1 = (z_{11}, \ldots, z_{1n})', \ldots, \boldsymbol{z}_m = (z_{m1}, \ldots, z_{mn})'$  for m > 1 under PIS. Consider the point estimates  $\boldsymbol{b}$ and RSS/(n-p), of  $\beta$  and  $\sigma^2$ , respectively, where  $\boldsymbol{b} = (\boldsymbol{X}'\boldsymbol{X})^{-1}\boldsymbol{X}'\boldsymbol{y}$  and RSS =  $(\boldsymbol{y} - \boldsymbol{X}\boldsymbol{b})'(\boldsymbol{y} - \boldsymbol{x})$  $Xb) = y'(I_n - P_X)y$  with  $I_k$  as the k-dimensional identity matrix and  $P_X = X(X'X)^{-1}X'$  is the orthogonal projection matrix to the column space of X. The synthetic data are obtained by drawing  $\boldsymbol{z}_1, \ldots, \boldsymbol{z}_m \stackrel{\text{iid}}{\sim} N_n\left(\boldsymbol{X}\boldsymbol{b}, \frac{\text{RSS}}{n-p}\boldsymbol{I}_n\right)$ . Equivalently, the synthetic data are obtained by drawing

 $z_{ji} \sim N(\boldsymbol{x}'_{i}\boldsymbol{b}, \frac{\text{RSS}}{n-p})$ , independently for i = 1, ..., n and j = 1, ..., m. Let  $\bar{z}_{i} = \frac{1}{m} \sum_{j=1}^{m} z_{ji}, S_{zi}^{2} = \sum_{j=1}^{m} (z_{ji} - \bar{z}_{i})^{2}$ , and  $S_{z}^{2} = \sum_{i=1}^{n} S_{zi}^{2}$ . If m > 1, then conditional on  $\boldsymbol{b}$  and RSS,

$$S_z^2 \sim \frac{\text{RSS}}{(n-p)} \chi^2_{n(m-1)}, \quad \bar{z}_i \sim N\left(\boldsymbol{x}'_i \boldsymbol{b}, \frac{\text{RSS}}{m(n-p)}\right), \ i = 1, \dots, n,$$

with these terms being (conditionally) independent. If m = 1, then  $\bar{z}_i = z_{1i}$  and  $S_{zi}^2 = 0$  for

 $i = 1, \dots, n, \text{ and hence } S_z^2 = 0.$   $\text{Let } \bar{\boldsymbol{z}} = (\bar{z}_1, \dots, \bar{z}_n)' \text{ and } \boldsymbol{b}_j^* = (\boldsymbol{X}'\boldsymbol{X})^{-1}\boldsymbol{X}'\boldsymbol{z}_j. \text{ We define } \overline{\boldsymbol{b}^*} = (\boldsymbol{X}'\boldsymbol{X})^{-1}\boldsymbol{X}'\bar{\boldsymbol{z}} = \frac{1}{m}\sum_{j=1}^m \boldsymbol{b}_j^* \text{ and } S_{\text{comb}}^2 = S_z^2 + m(\bar{\boldsymbol{z}} - \boldsymbol{X}\overline{\boldsymbol{b}^*})'(\bar{\boldsymbol{z}} - \boldsymbol{X}\overline{\boldsymbol{b}^*}), \text{ and conditional on } \boldsymbol{b} \text{ and RSS},$ 

$$\overline{\boldsymbol{b}^*} \sim N_p\left(\boldsymbol{b}, \frac{\mathrm{RSS}}{m(n-p)} (\boldsymbol{X}' \boldsymbol{X})^{-1}\right), \quad S_{\mathrm{comb}}^2 \sim \frac{\mathrm{RSS}}{(n-p)} \chi_{n(m-1)+n-p}^2$$

which are (conditionally) independent and are jointly sufficient for  $(\beta, \sigma^2)$ . From [2], we have the following result.

**Theorem 3.1.** The joint pdf of  $(\overline{b^*}, S_{\text{comb}}^2)$  is given by

$$f_{\boldsymbol{\beta},\sigma^{2}}(\overline{\boldsymbol{b}^{*}}, S_{\text{comb}}^{2}) \propto \int_{0}^{\infty} e^{-\frac{1}{2} \left[ \frac{(\overline{\boldsymbol{b}^{*}} - \boldsymbol{\beta})'(\boldsymbol{X}'\boldsymbol{X})(\overline{\boldsymbol{b}^{*}} - \boldsymbol{\beta})}{\sigma^{2}(1 + \frac{\psi}{m(n-p)})} + \frac{(n-p)S_{\text{comb}}^{2}}{\sigma^{2}\psi} + \psi \right]}{\frac{(S_{\text{comb}}^{2})^{\frac{nm-p}{2}-1}}{\sigma^{nm}\psi^{\frac{n(m-1)+p+2}{2}}}} \times \left[ 1 + \frac{m(n-p)}{\psi} \right]^{-p/2} d\psi$$

$$(1)$$

#### Posterior distributions of $\beta$ and $\sigma^2$ 3.1

For Bayesian inference on the other unknown parameters we assume non-informative improper priors and assume that all unknown quantities are a priori independent. Specifically, we assume

$$\pi(\boldsymbol{\beta}, \sigma^2) = \pi(\boldsymbol{\beta})\pi(\sigma^2)$$

where  $\pi(\beta) \propto 1$  and  $\pi(\sigma) \propto \sigma^{-\delta}$  and hence the induced prior on  $\sigma^2$  is  $\pi(\sigma^2) \propto (\sigma^2)^{-\frac{\delta+1}{2}}$  for  $\delta > 0$ . For doing posterior computation, we use latent variable augmentation. Consider the latent variable  $\psi = (\hat{\sigma}/\sigma)^2$  where  $\hat{\sigma}^2 = RSS/(n-p)$ . Then following the development in the single imputation case in [1], we use the joint distribution of  $(\overline{b^*}, S^2_{\text{comb}})$  conditional on the latent quantity  $\psi$  to perform posterior computation. Multiplying the prior with the conditional likelihood one obtains the posterior distributions as:

$$\boldsymbol{\beta} | \overline{\boldsymbol{b}^*}, S_{\text{comb}}^2, \sigma^2, \psi \sim N_p \left( \overline{\boldsymbol{b}^*}, \sigma^2 (1 + \frac{\psi}{m(n-p)}) (\boldsymbol{X}' \boldsymbol{X})^{-1} \right)$$

$$\sigma^2 | \overline{\boldsymbol{b}^*}, S_{\text{comb}}^2, \psi \sim \text{Scale-inv-} \chi^2 \left( nm - p + \delta - 1, \frac{(n-p)S_{\text{comb}}^2}{\psi(nm-p+\delta-1)} \right)$$

$$\psi \sim \chi_{n-p+\delta-1}^2$$
(2)

The posterior distributions are proper as long as  $n > \max\{p, p - \delta + 1\}$  (this also ensures that  $nm - p + \delta - 1 > 0$  since m > 1, which is necessary for the posterior distribution of  $\sigma^2$  to be proper). For m = 1 in the above formula yields the same results as obtained for the singly imputed plug-in sampling case in [1]. From the joint posterior of  $(\overline{b^*}, S^2_{\text{comb}}, \psi)$  the Bayes estimators of  $\beta$  and  $\sigma^2$  follows immediately:

$$\begin{aligned} \hat{\boldsymbol{\beta}}_{\text{BAYES}} &= \mathcal{E}(\boldsymbol{\beta} \,|\, \overline{\boldsymbol{b}^*}, S_{\text{comb}}^2) = \mathcal{E}_{\psi} \,\mathcal{E}_{\sigma^2} \,\mathcal{E}(\boldsymbol{\beta} \,|\, \overline{\boldsymbol{b}^*}, S_{\text{comb}}^2, \sigma^2, \psi) = \mathcal{E}_{\psi} \,\mathcal{E}_{\sigma^2}(\overline{\boldsymbol{b}^*}) = \overline{\boldsymbol{b}^*} \\ \hat{\sigma}_{\text{BAYES}}^2 &= \mathcal{E}(\sigma^2 \,|\, \overline{\boldsymbol{b}^*}, S_{\text{comb}}^2) = \mathcal{E}_{\psi} \,\mathcal{E}(\sigma^2 \,|\, \overline{\boldsymbol{b}^*}, S_{\text{comb}}^2, \psi) = \mathcal{E}_{\psi}(\frac{(n-p)S_{\text{comb}}^2}{\psi(nm-p+\delta-3)}) \\ &= \frac{(n-p)S_{\text{comb}}^2}{(nm-p+\delta-3)} \,\mathcal{E}_{\psi}(\frac{1}{\psi}) = \frac{(n-p)S_{\text{comb}}^2}{(nm-p+\delta-3)(n-p+\delta-3)} \end{aligned}$$

# **3.2** Credible Sets for $\beta$ and $\sigma^2$

In this section we describe how to get credible sets for the regression parameters. Let  $U \coloneqq \frac{1}{\sigma^2}$ . Then following (2), we have

$$U \mid S_{\text{comb}}^2, \psi \sim \Gamma\left(\frac{nm - p + \delta - 1}{2}, \frac{(n - p)S_{\text{comb}}^2}{2\psi}\right)$$
(3)

Define  $K \coloneqq \frac{S_{\text{comb}}^2}{\sigma^2} = U \times S_{\text{comb}}^2$ . Then from (3) and the fact that if  $X \sim \Gamma(\alpha, \beta)$  then  $cX \sim \Gamma(\alpha, \beta/c)$ , it follows that  $K \mid S_{\text{comb}}^2, \psi \sim \Gamma\left(\frac{nm-p+\delta-1}{2}, \frac{(n-p)}{2\psi}\right)$ . Since the right hand side is independent of  $S_{\text{comb}}^2$ , it follows that

$$K \mid \psi \sim \Gamma\left(\frac{nm - p + \delta - 1}{2}, \frac{(n - p)}{2\psi}\right).$$
(4)

Using (4) and the fact that  $\psi \sim \chi^2_{n-p+\delta-1}$ , an  $(1-\gamma)$  level credible set for  $\sigma^2$  based on K is

$$\left[\frac{S_{\text{comb}}^2}{b_{n,p,\delta;\gamma}}, \frac{S_{\text{comb}}^2}{a_{n,p,\delta;\gamma}}\right]$$

where  $a_{n,p,\delta;\gamma}$  and  $b_{n,p,\delta;\gamma}$  are any two constants that satisfy  $1 - \gamma = P(a_{n,p,\delta;\gamma} \leq K \leq b_{n,p,\delta;\gamma})$ . The length of the credible interval is  $S^2_{\text{comb}}\left(\frac{1}{a_{n,p,\delta;\gamma}} - \frac{1}{b_{n,p,\delta;\gamma}}\right)$ . To obtain a credible set for  $\beta$ , we define  $V \coloneqq \frac{(\beta - \overline{b^*})'(X'X)(\beta - \overline{b^*})}{\sigma^2(1 + \frac{\psi}{m(n-p)})}$  Then  $V \mid \overline{b^*}, S^2_{\text{comb}}, \sigma^2, \psi \sim \chi^2_p$  and thus unconditionally  $V \sim \chi^2_p$ . Also V is independent of  $(\overline{b^*}, S^2_{\text{comb}}, \sigma^2, \psi)$  and thus V is independent of U. If  $U^* \coloneqq \frac{U(n-p)S^2_{\text{comb}}}{\psi}$ then  $U^* | S^2_{\text{comb}}, \psi \sim \chi^2_{nm-p+\delta-1}$ , and unconditionally  $U^* \sim \chi^2_{nm-p+\delta-1}$ . As V is independent of U, it is independent of  $U^*$ . Finally we define the pivot for  $\beta$  as

$$T_m^2 \coloneqq \frac{(\boldsymbol{\beta} - \overline{\boldsymbol{b}^*})'(\boldsymbol{X}'\boldsymbol{X})(\boldsymbol{\beta} - \overline{\boldsymbol{b}^*})}{S_{\mathrm{comb}}^2}$$

. To derive the posterior distribution of  $T_m^2,$  note that, conditionally given  $\psi,$ 

$$T_m^2 = \frac{\sigma^2 V\left(1 + \frac{\psi}{m(n-p)}\right)}{\frac{\psi U^*}{U(n-p)}}$$
$$= \frac{V}{U^*} \left[\frac{1 + \frac{\psi}{m(n-p)}}{\psi}\right] (n-p)$$
$$\sim \frac{\chi_p^2}{\chi_{nm-p+\delta-1}^2} \left(\frac{1}{m} + \frac{n-p}{\psi}\right)$$
$$= \left[\frac{p}{nm-p+\delta-1}\right] F_{p,n-p+\delta-1} \left(\frac{1}{m} + \frac{n-p}{\psi}\right)$$

Hence the pivot for  $\beta$  is computed from the distribution of  $T_m^2$  which follows from

$$\psi \sim \chi^2_{n-p+\delta-1}$$
  
$$T_m^2 | \psi \sim \left[ \frac{p}{nm-p+\delta-1} \right] \left[ \frac{1}{m} + \frac{n-p}{\psi} \right] F_{p,nm-p+\delta-1}.$$

A  $(1 - \gamma)$  level credible ellipsoid for  $\beta$  based on  $T_m^2$  is given by

$$\{\boldsymbol{\beta} : T_m^2 \leq d_{n,p,\delta,m;\gamma}\}$$

where  $d_{n,p,\delta,m;\gamma}$  satisfies  $1 - \gamma = P(T_m^2 \leq d_{n,p,\delta,m;\gamma})$ . The volume of the credible ellipsoid is

$$V_{\boldsymbol{\beta}}(\boldsymbol{z}_1,\ldots,\boldsymbol{z}_m,\boldsymbol{X}) = \frac{\pi^{p/2}}{\Gamma\left(\frac{p}{2}+1\right)} \left(d_{n,p,\delta,m;\gamma} S_{\text{comb}}^2\right)^{p/2} \left|\boldsymbol{X}'\boldsymbol{X}\right|^{-1/2}$$

**Remark 3.1.** If one is interested in the credible set of a single regression coefficient or more generally in the credible set of a linear combination of  $\beta$ , namely,  $A\beta = \eta$  where A is a  $k \times p$ dimensional matrix with rank(A) = k < p, we define  $T_{m,\eta}^2 = (\eta - A\overline{b^*})' \{A(XX')^{-1}A'\}^{-1}(\eta - A\overline{b^*})/S_{comb}^2$ , and proceed by noting that

$$T_{m,\boldsymbol{\eta}}^2 | \psi \sim \left[ \frac{k}{nm - p + \delta - 1} \right] \left[ \frac{1}{m} + \frac{n - p}{\psi} \right] F_{k,nm - p + \delta - 1} \quad and \quad \psi \sim \chi_{n - p + \delta - 1}^2$$

## 4 Posterior Predictive Sampling method

In this section we repeat the analysis done under PIS for PPS. We consider the setup described in Section 3. The synthetic data are generated by repeating the following steps below independently for each j = 1, ..., m.

(a) Draw  $(\boldsymbol{\beta}_{i}^{*}, \sigma_{i}^{*2})$  from the posterior distribution (2).

(b) Draw 
$$\boldsymbol{z}_j = (z_{j1}, \dots, z_{jn})' \sim N_n(\boldsymbol{X}\boldsymbol{\beta}_j^*, \sigma_j^{*2}\boldsymbol{I}_n).$$

The released synthetic data are  $\mathbf{z}_1, \ldots, \mathbf{z}_m$  along with the matrix of predictor variables  $\mathbf{X}$ . The sufficient statistics for the synthetic data are:  $\mathbf{b}_j^* = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{z}_j$  and  $\mathrm{RSS}_j^* = (\mathbf{z}_j - \mathbf{X}\mathbf{b}_j^*)'(\mathbf{z}_j - \mathbf{X}\mathbf{b}_j^*)$ , for  $j = 1, \ldots, m$ . It can be shown that  $(\mathbf{b}_1^*, \mathrm{RSS}_1^*), \ldots, (\mathbf{b}_m^*, \mathrm{RSS}_m^*)$  are jointly sufficient for  $(\boldsymbol{\beta}, \sigma^2)$ . In view of the sampling mechanism above, the joint distribution of  $\mathbf{b}_1^*, \ldots, \mathbf{b}_m^*, \mathrm{RSS}_1^*, \ldots, \mathrm{RSS}_m^*$ ,  $\boldsymbol{\beta}_1^*, \ldots, \boldsymbol{\beta}_m^*, \sigma_1^{*2}, \ldots, \sigma_m^{*2}, \mathbf{b}$  and RSS has the following hierarchical structure:

$$\begin{split} \boldsymbol{b}_{j}^{*} | \operatorname{RSS}_{1}^{*}, \dots, \operatorname{RSS}_{m}^{*}, \boldsymbol{\beta}_{1}^{*}, \dots, \boldsymbol{\beta}_{m}^{*}, \sigma_{1}^{*2}, \dots, \sigma_{m}^{*2}, \boldsymbol{b}, \operatorname{RSS} \ \sim \ N_{p}(\boldsymbol{\beta}_{j}^{*}, \sigma_{j}^{*2}(\boldsymbol{X}'\boldsymbol{X})^{-1}), \\ \operatorname{RSS}_{j}^{*} | \boldsymbol{\beta}_{1}^{*}, \dots, \boldsymbol{\beta}_{m}^{*}, \sigma_{1}^{*2}, \dots, \sigma_{m}^{*2}, \boldsymbol{b}, \operatorname{RSS} \ \sim \ \sigma_{j}^{*2} \chi_{n-p}^{2}, \\ \boldsymbol{\beta}_{j}^{*} | \sigma_{1}^{*2}, \dots, \sigma_{m}^{*2}, \boldsymbol{b}, \operatorname{RSS} \ \sim \ N_{p}(\boldsymbol{b}, \sigma_{j}^{*2}(\boldsymbol{X}'\boldsymbol{X})^{-1}), \\ \sigma_{j}^{*2} | \boldsymbol{b}, \operatorname{RSS} \ \sim \ \frac{\operatorname{RSS}}{\chi_{n-p+\alpha-1}^{2}}, \\ \boldsymbol{b} \ \sim \ N_{p}(\boldsymbol{\beta}, \sigma^{2}(\boldsymbol{X}'\boldsymbol{X})^{-1}), \\ \operatorname{RSS} \ \sim \ \sigma^{2} \chi_{n-p}^{2}. \end{split}$$

which are generated independently for  $j = 1, \ldots, m$ , whenever applicable. Hence,

$$\begin{split} f(\boldsymbol{b}_{1}^{*},\dots,\boldsymbol{b}_{m}^{*},\mathrm{RSS}_{1}^{*},\dots,\mathrm{RSS}_{m}^{*},\boldsymbol{\beta}_{1}^{*},\dots,\boldsymbol{\beta}_{m}^{*},\sigma_{1}^{*2},\dots,\sigma_{m}^{*2},\boldsymbol{b},\mathrm{RSS}) &= \\ & \prod_{j=1}^{m} (2\pi\sigma_{j}^{*2})^{-p/2} |\boldsymbol{X}'\boldsymbol{X}|^{1/2} \exp\left[-\frac{1}{2\sigma_{j}^{*2}}(\boldsymbol{b}_{j}^{*}-\boldsymbol{\beta}_{j}^{*})'(\boldsymbol{X}'\boldsymbol{X})(\boldsymbol{b}_{j}^{*}-\boldsymbol{\beta}_{j}^{*})\right] \\ & \times \prod_{j=1}^{m} \frac{(\mathrm{RSS}_{j}^{*})^{\frac{n-p}{2}-1}}{2^{\frac{n-p}{2}} \Gamma(\frac{n-p}{2})} (\sigma_{j}^{*2})^{-(n-p)/2} \exp\left[-\frac{\mathrm{RSS}_{j}^{*}}{2\sigma_{j}^{*2}}\right] \\ & \times \prod_{j=1}^{m} (2\pi\sigma_{j}^{*2})^{-p/2} |\boldsymbol{X}'\boldsymbol{X}|^{1/2} \exp\left[-\frac{1}{2\sigma_{j}^{*2}}(\boldsymbol{\beta}_{j}^{*}-\boldsymbol{b})'(\boldsymbol{X}'\boldsymbol{X})(\boldsymbol{\beta}_{j}^{*}-\boldsymbol{b})\right] \\ & \times \prod_{j=1}^{m} \frac{(\mathrm{RSS})^{(n-p+\alpha-1)/2}}{2^{(n-p+\alpha-1)/2} \Gamma(\frac{n-p+\alpha-1}{2})} (\sigma_{j}^{*2})^{-(n-p+\alpha-1)/2-1} \exp\left[-\frac{\mathrm{RSS}}{2\sigma_{j}^{*2}}\right] \\ & \times (2\pi\sigma^{2})^{-p/2} |\boldsymbol{X}'\boldsymbol{X}|^{1/2} \exp\left[-\frac{1}{2\sigma^{2}}(\boldsymbol{b}-\boldsymbol{\beta})'(\boldsymbol{X}'\boldsymbol{X})(\boldsymbol{b}-\boldsymbol{\beta})\right] \\ & \times \frac{(\mathrm{RSS})^{\frac{n-p}{2}-1}}{2^{\frac{n-p}{2}} \Gamma(\frac{n-p}{2})} (\sigma^{2})^{-(n-p)/2} \exp\left[-\frac{\mathrm{RSS}}{2\sigma^{2}}\right] \end{split}$$

Patterned after Theorem 3.1, integrating out  $\beta_j^*$ 's, **b**, RSS, we get the relevant likelihood as

$$\begin{split} L(\boldsymbol{\beta}, \sigma^{2}, \sigma_{j}^{*2}, j = 1, ..., m \mid \mathbf{b}_{j}^{*}, \mathrm{RSS}_{j}^{*}, j = 1, ..., m) &= \\ & (2\pi)^{-p/2} \frac{\left(\frac{1}{\sigma^{2}}\right)^{p/2} \left(\sum_{j=1}^{m} \frac{1}{2\sigma_{j}^{*2}}\right)^{p/2}}{\left(\frac{1}{\sigma^{2}} + \sum_{j=1}^{m} \frac{1}{2\sigma_{j}^{*2}}\right)^{p/2}} \mid \mathbf{X}' \mathbf{X} \mid^{1/2} \\ & \times \exp\left[-\frac{1}{2} \frac{\left(\frac{1}{\sigma^{2}}\right) \left(\sum_{j=1}^{m} \frac{1}{2\sigma_{j}^{*2}}\right)}{\left(\frac{1}{\sigma^{2}} + \sum_{j=1}^{m} \frac{1}{2\sigma_{j}^{*2}}\right)} \left(\boldsymbol{\beta} - \frac{\sum_{j=1}^{m} \frac{\mathbf{b}_{j}^{*}}{2\sigma_{j}^{*2}}}{\sum_{j=1}^{m} \frac{1}{2\sigma_{j}^{*2}}}\right)' (\mathbf{X}' \mathbf{X}) \left(\boldsymbol{\beta} - \frac{\sum_{j=1}^{m} \frac{\mathbf{b}_{j}^{*}}{2\sigma_{j}^{*2}}}{\sum_{j=1}^{m} \frac{1}{2\sigma_{j}^{*2}}}\right)\right] \\ & \times \prod_{j=1}^{m} (2\pi)^{-p/2} \left(2\sigma_{j}^{*2}\right)^{-p/2} \mid \mathbf{X}' \mathbf{X} \mid^{1/2} \exp\left[-\frac{1}{2} \cdot \frac{1}{2\sigma_{j}^{*2}} \left(\mathbf{b}_{j}^{*} - \frac{\sum_{j=1}^{m} \frac{\mathbf{b}_{j}^{*}}{2\sigma_{j}^{*2}}}{\sum_{j=1}^{m} \frac{1}{2\sigma_{j}^{*2}}}\right)' (\mathbf{X}' \mathbf{X}) \left(\mathbf{b}_{j}^{*} - \frac{\sum_{j=1}^{m} \frac{\mathbf{b}_{j}^{*}}{2\sigma_{j}^{*2}}}{\sum_{j=1}^{m} \frac{1}{2\sigma_{j}^{*2}}}\right)\right] \\ & \times \frac{\left(\sum_{j=1}^{m} \frac{1}{2\sigma_{j}^{*2}}\right)^{-\frac{p}{2}} \left(\frac{1}{\sigma^{2}} + \sum_{j=1}^{m} \frac{1}{\sigma_{j}^{*2}}\right)^{-\frac{(m+1)(n-p)+m(\alpha-1)}{2}} \exp\left[-\sum_{j=1}^{m} \frac{\mathrm{RSS}_{j}^{*}}{2\sigma_{j}^{*2}}\right] \\ & \times \frac{\Gamma\left(\frac{m(n-p+\alpha-1)+(n-p)}{2}\right)}{\left(\Gamma\left(\frac{n-p+\alpha-1}{2}\right)\right)^{m} \left(\Gamma\left(\frac{n-p}{2}\right)\right)^{m+1} (2\pi)^{p/2} \mid \mathbf{X}' \mathbf{X} \mid^{-1/2} \left(\prod_{j=1}^{m} \frac{(\mathrm{RSS}_{j}^{*})^{\frac{n-p}{2}-1}}{2^{\frac{n-p}{2}}}\right). \tag{5}$$

Observe that the quantity inside the exponential in the second line vanishes for m = 1. Multiplying the likelihood in (5) by the prior  $\pi(\beta, \sigma^2) \propto (\sigma^2)^{-\frac{\delta+1}{2}}$ , and separating the parameters we get the posteriors as:

$$\beta \mid \sigma^{2}, \sum_{j=1}^{m} \frac{\boldsymbol{b}_{j}^{*}}{\sigma_{j}^{*2}}, \sum_{j=1}^{m} \frac{1}{\sigma_{j}^{*2}} \sim N_{p} \left( \frac{\sum_{j=1}^{m} \frac{\boldsymbol{b}_{j}^{*}}{\sigma_{j}^{*2}}}{\sum_{j=1}^{m} \frac{1}{\sigma_{j}^{*2}}}, \left( 1 + \sigma^{2} \left( \sum_{j=1}^{m} \frac{1}{\sigma_{j}^{*2}} \right) \right) (\boldsymbol{X}'\boldsymbol{X})^{-1} \right),$$

$$\sigma^{2} \left( \sum_{j=1}^{m} \frac{1}{\sigma_{j}^{*2}} \right) \mid \left( \sum_{j=1}^{m} \frac{1}{\sigma_{j}^{*2}} \right) \sim \beta' \left( \frac{m(n-p+\alpha-1)-\delta+1}{2}, \frac{n-p+\delta-1}{2} \right)$$
(6)

Thus,  $\sigma^2 \left( \sum_{j=1}^m 1/\sigma_j^{*2} \right)$  is independent of both latent variables and data. The posterior distributions are proper as long as  $n > \max\left\{p, p - \delta + 1, p - \alpha + 1, p - \alpha + 1 + \frac{\delta - 1}{m}\right\}$ . The latent variables

have the following distribution:

$$g(\sigma_{1}^{*2}, \dots, \sigma_{m}^{*2} | \boldsymbol{b}_{1}^{*}, \dots, \boldsymbol{b}_{m}^{*}, \operatorname{RSS}_{1}^{*}, \dots, \operatorname{RSS}_{m}^{*}) \propto \frac{\left(\prod_{j=1}^{m} (\sigma_{j}^{*2})\right)^{-\frac{2n-p+\alpha+1}{2}}}{\left(\sum_{j=1}^{m} \frac{1}{\sigma_{j}^{*2}}\right)^{\frac{m(n-p+\alpha-1)+p-\delta}{2}}} \exp\left[-\frac{1}{4}\sum_{j=1}^{m} \frac{1}{\sigma_{j}^{*2}} \left\{ \left(\boldsymbol{b}_{j}^{*} - \frac{\sum_{j=1}^{m} \frac{\boldsymbol{b}_{j}^{*}}{\sigma_{j}^{*2}}}{\sum_{j=1}^{m} \frac{1}{\sigma_{j}^{*2}}}\right)'(\boldsymbol{X}\boldsymbol{X}') \left(\boldsymbol{b}_{j}^{*} - \frac{\sum_{j=1}^{m} \frac{\boldsymbol{b}_{j}^{*}}{\sigma_{j}^{*2}}}{\sum_{j=1}^{m} \frac{1}{\sigma_{j}^{*2}}}\right) + 2\operatorname{RSS}_{j}^{*}\right\}\right].$$

Defining  $\sigma_j^2 = 1/\sigma_j^{*2}$ , we have the joint posterior of the transformed latent variables as

$$h(\sigma_1^2, \dots, \sigma_m^2 | \boldsymbol{b}_1^*, \dots, \boldsymbol{b}_m^*, \operatorname{RSS}_1^*, \dots, \operatorname{RSS}_m^*) \propto \frac{\left(\prod_{j=1}^m (\sigma_j^2)\right)^{\frac{2n-p+\alpha-1}{2}}}{\left(\sum_{j=1}^m \sigma_j^2\right)^{\frac{m(n-p+\alpha-1)+p-\delta}{2}}} \exp\left[-\frac{1}{4}\sum_{j=1}^m \sigma_j^2 \left\{ \left(\boldsymbol{b}_j^* - \frac{\sum_{j=1}^m \sigma_j^2 \boldsymbol{b}_j^*}{\sum_{j=1}^m \sigma_j^2}\right)' (\boldsymbol{X}\boldsymbol{X}') \left(\boldsymbol{b}_j^* - \frac{\sum_{j=1}^m \sigma_j^2 \boldsymbol{b}_j^*}{\sum_{j=1}^m \sigma_j^2}\right) + 2\operatorname{RSS}_j^* \right\}\right].$$
(7)

Let us denote the quantity inside the exponential of (7) as Q. Our goal here is to sample from (7).

#### 4.1 Approach I:

To suitably transform  $\sigma_1^2, \ldots, \sigma_m^2$  to draw a sample, our first method is to use

$$v_1 = \sigma_1^2, v_2 = \frac{\sigma_2^2}{\sigma_1^2}, \dots, v_m = \frac{\sigma_m^2}{\sigma_1^2}$$
 (8)

so that  $(v_1, v_2, \ldots, v_m)$  can be sampled as

$$v_{1} | v_{2}, \dots, v_{m}, \text{data} \sim \frac{2\chi_{nm-p+\delta-1}^{2}}{Q},$$
  

$$\pi(v_{2}, \dots, v_{m} | \text{data}) \leq \frac{\mathbf{B}\left(\frac{2n-p+\alpha-1}{2}, \dots, \frac{2n-p+\alpha-1}{2}\right)}{(\text{RSS}_{\min}^{*})^{\frac{nm-p+\delta-1}{2}}} g_{\text{InvDir}}(v_{2}, \dots, v_{m}),$$
(9)

where  $g_{\text{InvDir}}(v_2, \ldots, v_m)$  is the pdf of an  $m^{\text{th}}$  order Inverse-Dirichlet  $\left(\frac{2n-p+\alpha-1}{2}, \ldots, \frac{2n-p+\alpha-1}{2}\right)$  distribution.

### 4.2 Approach II:

One could also use the following transformation:

$$v_1 = \sigma_1^2, v_2 = \frac{\sigma_1^2}{\sigma_2^2}, v_3 = \frac{\sigma_1^2 + \sigma_2^2}{\sigma_3^2}, \dots, v_m = \frac{\sigma_1^2 + \dots + \sigma_{m-1}^2}{\sigma_m^2},$$
(10)

so that  $(v_1, v_2, \ldots, v_m)$  can be sampled as

$$v_{1} | v_{2}, \dots, v_{m}, \text{data} \sim \frac{2 \chi_{nm-p+\delta-1}^{2}}{Q},$$
  

$$\pi(v_{2}, \dots, v_{m} | \text{data}) \leq \frac{\prod_{j=2}^{m} B\left(\frac{(2n-p+\alpha-1)(j-1)}{2}, \frac{2n-p+\alpha-1}{2}\right)}{(\text{RSS}_{\min}^{*})^{\frac{nm-p+\delta-1}{2}}} \left(\prod_{j=2}^{m} g_{j}'(v_{j})\right), \quad (11)$$

where  $g'_j(v_j)$  is the pdf of a Beta-Prime  $\left(\frac{(2n-p+\alpha-1)(j-1)}{2}, \frac{2n-p+\alpha-1}{2}\right)$  distribution, independently for  $j = 2, \ldots, m$ . Since if  $X \sim \beta'(a, b)$  then  $X^{-1} \sim \beta'(b, a)$ , the reciprocal transformation would also work in (11).

## 5 Partially Sensitive Data

In this section we discuss the situation when only a part of the response vector  $\mathbf{y}$ , say  $(y_1, \ldots, y_r)$  is sensitive and hence needs to be protected. We provide two methods, depending on the nature of imputation of the synthetic data.

#### Method I: Using whole data estimates to impute synthetic data

#### Plug-In Sampling

The original data now has the same setup as in Section 3 with both r > p and n - r > p assumed to hold. We synthesize m copies of the original data  $\boldsymbol{y} = (\boldsymbol{y}_1, \boldsymbol{y}_2)$  given by  $\{\boldsymbol{y}^{*j} = (\boldsymbol{y}_1^{*j}, \boldsymbol{y}_2) : j = 1, \ldots, m\}$  whose sufficient statistics are given by  $(\boldsymbol{b}_1^{*1}, \ldots, \boldsymbol{b}_1^{*m}, \text{RSS}_1^{*1}, \ldots, \text{RSS}_1^{*m}, \boldsymbol{b}_2, \text{RSS}_2)$ . We denote  $\overline{\boldsymbol{b}_1^*} = \frac{1}{m} \sum_{j=1}^m \boldsymbol{b}_1^{*j}$ . Then we can derive the following posterior distributions in a similar manner as before

$$\begin{aligned} \boldsymbol{\beta} \,|\, \boldsymbol{\sigma}^{2}, \psi, \overline{\boldsymbol{b}_{1}^{*}}, \boldsymbol{b}_{2} &\sim N_{p} \left[ \left( \frac{\boldsymbol{X}_{1}^{\prime} \boldsymbol{X}_{1}}{1+\psi} + \frac{\boldsymbol{X}_{2}^{\prime} \boldsymbol{X}_{2}}{m} \right)^{-1} \left( \frac{\boldsymbol{X}_{1}^{\prime} \boldsymbol{X}_{1}}{1+\psi} \overline{\boldsymbol{b}_{1}^{*}} + \frac{\boldsymbol{X}_{2}^{\prime} \boldsymbol{X}_{2}}{m} \boldsymbol{b}_{2} \right), \frac{\boldsymbol{\sigma}^{2}}{m} \left( \frac{\boldsymbol{X}_{1}^{\prime} \boldsymbol{X}_{1}}{1+\psi} + \frac{\boldsymbol{X}_{2}^{\prime} \boldsymbol{X}_{2}}{m} \right)^{-1} \right], \\ \boldsymbol{\sigma}^{2} \,|\, \psi, \boldsymbol{b}_{1}^{*1}, \dots, \boldsymbol{b}_{1}^{*m}, \text{RSS}_{1}^{*1}, \dots, \text{RSS}_{1}^{*m}, \boldsymbol{b}_{2}, \text{RSS}_{2} &\sim \text{Scale-inv-} \boldsymbol{\chi}^{2} \left( \nu, \tau_{1}^{2} \right), \\ \boldsymbol{\pi}(\psi \,|\, \boldsymbol{b}_{1}^{*1}, \dots, \boldsymbol{b}_{1}^{*m}, \text{RSS}_{1}^{*1}, \dots, \text{RSS}_{1}^{*m}, \boldsymbol{b}_{2}, \text{RSS}_{2}), \\ \propto \left| \frac{\boldsymbol{X}_{1}^{\prime} \boldsymbol{X}_{1}}{1+\psi} + \frac{\boldsymbol{X}_{2}^{\prime} \boldsymbol{X}_{2}}{m} \right|^{-\frac{1}{2}} \psi^{-\frac{(m-1)(r-p)}{2}-1} (1+\psi)^{-\frac{mp}{2}} e^{\frac{(r-p)\psi}{2}} \left\{ \nu \tau_{1}^{2} \right\}^{-\frac{\nu}{2}} \end{aligned}$$

where  $\nu = n + (m-1)r - p + \delta - 1$  and

$$\nu \tau_1^2 = \sum_{j=1}^m \left( \boldsymbol{b}_1^{*j} - \overline{\boldsymbol{b}_1^*} \right)' \frac{\boldsymbol{X}_1' \boldsymbol{X}_1}{1 + \psi} \left( \boldsymbol{b}_1^{*j} - \overline{\boldsymbol{b}_1^*} \right) + m \left( \overline{\boldsymbol{b}_1^*} - \boldsymbol{b}_2 \right)' \left( (1 + \psi) \left( \boldsymbol{X}_1' \boldsymbol{X}_1 \right)^{-1} + \left( \boldsymbol{X}_2' \boldsymbol{X}_2 \right)^{-1} \right)^{-1} \left( \boldsymbol{b}_1^{*j} - \overline{\boldsymbol{b}_1^*} \right) + \sum_{j=1}^m \frac{\text{RSS}_1^{*j}}{\psi} + \text{RSS}_2.$$

The posterior distributions are proper as long as r > p,  $n - r > \max\{p, p - rm - \delta + 1\}$ .

Let 
$$\mathbf{X}'\mathbf{X}_{\psi,m} = \frac{\mathbf{X}_1'\mathbf{X}_1}{1+\psi} + \frac{\mathbf{X}_2'\mathbf{X}_2}{m}$$
. Note that  
 $\mathbf{X}'\mathbf{X}_{\psi,m} \ge \frac{\mathbf{X}'\mathbf{X}}{m(1+\psi)} \implies |\mathbf{X}'\mathbf{X}_{\psi,m}|^{-\frac{1}{2}} \le m^{\frac{p}{2}}(1+\psi)^{\frac{p}{2}} |\mathbf{X}'\mathbf{X}|^{-\frac{1}{2}} \nu \tau_1^2 \ge \sum_{j=1}^m \frac{\mathrm{RSS}_1^{*j}}{\psi}$ 

. Here  $A \ge B$  is used to mean (A - B) is positive semidefinite. Now using the fact that  $(1 + A) \ge B$  $\psi$ )<sup> $-\frac{(m-1)p}{2}$ </sup>  $\leq 1$  (as  $\psi > 0$ ), we can sample from the distribution of latent variables using the Accept-Reject algorithm as follows:

$$\pi(\psi \,|\, \text{data}) \leq \frac{m^{\frac{p}{2}} \,|\boldsymbol{X}'\boldsymbol{X}|^{-\frac{1}{2}} \,2^{\frac{n+(m-2)p+\delta-1}{2}} \Gamma\left(\frac{n+(m-2)p+\delta-1}{2}\right)}{(r-p)^{\frac{n+(m-2)p+\delta-1}{2}} \left(\sum_{j=1}^{m} \text{RSS}_{1}^{*j}\right)^{\frac{n+(m-1)r-p+\delta-1}{2}} \tilde{f}_{\text{ScaledChi}}(\psi)$$

where  $\tilde{f}_{\text{ScaledChi}}(\psi)$  is the pdf of a Gamma  $\left(\frac{n+(m-2)p+\delta-1}{2}, \frac{r-p}{2}\right)$  distribution. We also assume  $n+(m-2)p+\delta-1>0$ . All expressions in the partially sensitive case coincide with the results in the fully sensitive case when all of y is sensitive.

#### **Posterior Predictive Sampling**

We follow the same process as in Section 5 for the PPS case to derive the following posterior distributions

$$\begin{split} \boldsymbol{\beta} \, | \, \sigma^2, \psi_1, \dots, \psi_m, \boldsymbol{b}_1^{*1}, \dots, \boldsymbol{b}_1^{*m}, \boldsymbol{b}_2 \\ \sim N_p \left[ \left( \sum_{j=1}^m \frac{\boldsymbol{X}_1' \boldsymbol{X}_1}{1+2\psi_j} + \boldsymbol{X}_2' \boldsymbol{X}_2 \right)^{-1} \left( \sum_{j=1}^m \frac{\boldsymbol{X}_1' \boldsymbol{X}_1}{1+2\psi_j} \boldsymbol{b}_1^{*j} + \boldsymbol{X}_2' \boldsymbol{X}_2 \boldsymbol{b}_2 \right), \sigma^2 \left( \sum_{j=1}^m \frac{\boldsymbol{X}_1' \boldsymbol{X}_1}{1+2\psi_j} + \boldsymbol{X}_2' \boldsymbol{X}_2 \right)^{-1} \right], \\ \sigma^2 \, | \, \psi_1, \dots, \psi_m, \boldsymbol{b}_1^{*1}, \dots, \boldsymbol{b}_1^{*m}, \operatorname{RSS}_1^{*1}, \dots, \operatorname{RSS}_1^{*m}, \boldsymbol{b}_2, \operatorname{RSS}_2 \sim \operatorname{Scale-inv-} \chi^2 \left( \nu, \tau_1^2 \right), \\ \pi(\psi_1, \dots, \psi_m \, | \, \boldsymbol{b}_1^{*1}, \dots, \boldsymbol{b}_1^{*m}, \operatorname{RSS}_1^{*1}, \dots, \operatorname{RSS}_1^{*m}, \boldsymbol{b}_2, \operatorname{RSS}_2) \\ \propto \left| \sum_{j=1}^m \frac{\boldsymbol{X}_1' \boldsymbol{X}_1}{1+2\psi_j} + \boldsymbol{X}_2' \boldsymbol{X}_2 \right|^{-\frac{1}{2}} \left( \prod_{j=1}^m \psi_j^{-1} (1+2\psi_j)^{-\frac{p}{2}} (1+\psi_j)^{-\frac{2r-2p+\alpha-1}{2}} \right) \left\{ \nu \tau_2^2 \right\}^{-\frac{\nu}{2}} \\ \text{where } \nu = n + (m-1)r - p + \delta - 1 \text{ and} \end{split}$$

$$\nu \tau_2^2$$

$$= \sum_{j=1}^m \left( \boldsymbol{b}_1^{*j} - \left( \sum_{j=1}^m \frac{1}{1+2\psi_j} \right)^{-1} \left( \sum_{j=1}^m \frac{\boldsymbol{b}_1^{*j}}{1+2\psi_j} \right) \right)' \frac{\boldsymbol{X}_1' \boldsymbol{X}_1}{1+2\psi_j} \left( \boldsymbol{b}_1^{*j} - \left( \sum_{j=1}^m \frac{1}{1+2\psi_j} \right)^{-1} \left( \sum_{j=1}^m \frac{\boldsymbol{b}_1^{*j}}{1+2\psi_j} \right) \right)$$

$$+ \left( \left( \left( \sum_{j=1}^m \frac{1}{1+2\psi_j} \right)^{-1} \left( \sum_{j=1}^m \frac{\boldsymbol{b}_1^{*j}}{1+2\psi_j} \right) - \boldsymbol{b}_2 \right)' \left( \left( \sum_{j=1}^m \frac{1}{1+2\psi_j} \right)^{-1} \left( \boldsymbol{X}_1' \boldsymbol{X}_1 \right)^{-1} + \left( \boldsymbol{X}_2' \boldsymbol{X}_2 \right)^{-1} \right)^{-1} \right)$$

$$\left( \left( \left( \sum_{j=1}^m \frac{1}{1+2\psi_j} \right)^{-1} \left( \sum_{j=1}^m \frac{\boldsymbol{b}_1^{*j}}{1+2\psi_j} \right) - \boldsymbol{b}_2 \right) + \sum_{j=1}^m \frac{\mathrm{RSS}_1^{*j}}{\psi_j} + \mathrm{RSS}_2.$$

The posterior distributions are proper as long as  $r > \max\left\{p, p - \alpha + 1, \frac{n + (2m-1)p - \alpha m + \delta + m - 1}{m+1}\right\}, n - r > \max\{p, p - rm - \delta + 1\}$ . Note that

$$\left(\sum_{j=1}^m \frac{1}{1+2\psi_j}\right) \boldsymbol{X}_1' \boldsymbol{X}_1 + \boldsymbol{X}_2' \boldsymbol{X}_2 \ge \left(\prod_{j=1}^m \frac{1}{1+2\psi_j}\right) \boldsymbol{X}_1' \boldsymbol{X}_1 + \boldsymbol{X}_2' \boldsymbol{X}_2 \ge \left(\prod_{j=1}^m \frac{1}{1+2\psi_j}\right) \boldsymbol{X}' \boldsymbol{X}_1 + \boldsymbol{X}_2' \boldsymbol{X}_2 \ge \left(\prod_{j=1}^m \frac{1}{1+2\psi_j}\right) \boldsymbol{X}_1' \boldsymbol{X}_2 + \boldsymbol{X}_2' \boldsymbol{X}_2 \ge \left(\prod_{j=1}^m \frac{1}{1+2\psi_j}\right) \boldsymbol{X}_2' \boldsymbol{X}_2 + \boldsymbol{X}_2' \boldsymbol{X}_2 \ge \left(\prod_{j=1}^m \frac{1}{1+2\psi_j}\right) \boldsymbol{X}_2' \boldsymbol{X}_2 = \left(\prod_{j=1}^m \frac{1}{1+2\psi_j}\right) \boldsymbol{X}_2' \boldsymbol{X}_2' \boldsymbol{X}_2 = \left(\prod_{j=1}^m \frac{1}{1+2\psi_j}\right) \boldsymbol{X}_2' \boldsymbol{X}_2 = \left(\prod_{j=1}^m \frac{1}{1+2\psi_j}\right) \boldsymbol{X}_2' \boldsymbol{X}_2'$$

and hence using rejection sampling one can sample the latent variables as:

$$\pi(\psi_1, \dots, \psi_m \,|\, \text{data}) \le \frac{\left(\prod_{j=1}^m B\left(\frac{n + (m-1)r - p + \delta - 1}{2m}, \frac{(m+1)r - n - (2m-1)p + m\alpha - \delta - m + 1}{2m}\right)\right)}{|\mathbf{X}' \mathbf{X}|^{\frac{1}{2}} \left(m \prod_{j=1}^m \text{RSS}_j^{*\frac{1}{m}}\right)^{\frac{n + (m-1)r - p + \delta - 1}{2}} \left(\prod_{j=1}^m \overline{g_j}'(\psi_j)\right),$$

where  $\overline{g_j}'(\psi_j)$  is the pdf of a Beta-Prime  $\left(\frac{n+(m-1)r-p+\delta-1}{2m}, \frac{(m+1)r-n-(2m-1)p+m\alpha-\delta-m+1}{2m}\right)$  distribution, independently for  $j = 1, \ldots, m$ .

### Method II: Using only estimates of sensitive part to impute synthetic data

#### Plug-In Sampling

As in Section 5, our analysis in this case will be based solely on the synthetic part. We require only r > p. The posterior distributions are given by

$$\begin{split} \boldsymbol{\beta} \, | \, \boldsymbol{\sigma}^2, \boldsymbol{\psi}, \overline{\boldsymbol{b}_1^*} &\sim & N_p \left( \overline{\boldsymbol{b}_1^*}, \boldsymbol{\sigma}^2 \left( (\boldsymbol{X}' \boldsymbol{X})^{-1} + \frac{\boldsymbol{\psi}}{m} (\boldsymbol{X}_1' \boldsymbol{X}_1)^{-1} \right) \right), \\ \boldsymbol{\sigma}^2 \, | \, \boldsymbol{\psi}, \boldsymbol{b}_1^{*1}, \dots, \boldsymbol{b}_1^{*m}, \operatorname{RSS}_1^{*1}, \dots, \operatorname{RSS}_1^{*m} &\sim \quad \operatorname{Scale-inv} \cdot \boldsymbol{\chi}^2 \left( \tilde{\boldsymbol{\nu}}, \tilde{\tau}_1^2 \right), \\ \boldsymbol{\psi} &\sim \quad \frac{\boldsymbol{\chi}_{n-p+\delta-1}^2}{n-p} \equiv \Gamma \left( \frac{n-p+\delta-1}{2}, \frac{n-p}{2} \right) \\ \text{where } \tilde{\boldsymbol{\nu}} = rm - p + \delta - 1, \, \tilde{\boldsymbol{\nu}} \tilde{\tau_1}^2 = \frac{1}{\psi} \left( \sum_{j=1}^m \operatorname{RSS}_1^{*j} + \sum_{j=1}^m \left( \boldsymbol{b}_1^{*j} - \overline{\boldsymbol{b}_1^*} \right)' (\boldsymbol{X}_1' \boldsymbol{X}_1) \left( \boldsymbol{b}_1^{*j} - \overline{\boldsymbol{b}_1^*} \right) \right). \text{ The } \end{split}$$

posterior distributions are proper as long as  $r > \max\left\{p, \frac{p-\delta+1}{m}\right\}, n > p-\delta+1$  and the results match our expressions from Section 3 when r = n.

#### **Posterior Predictive Sampling**

Following the derivations in the plug-in sampling case for partially sensitive response, one can derive analogous expressions for the posterior in the posterior predictive sampling scheme. The posteriors are

$$\beta | \sigma^{2}, \psi_{1}, \dots, \psi_{m}, \boldsymbol{b}_{1}^{*1}, \dots, \boldsymbol{b}_{1}^{*m}$$

$$\sim N_{p} \left[ \left( \sum_{j=1}^{m} \psi_{j}^{-1} \right)^{-1} \left( \sum_{j=1}^{m} \frac{\boldsymbol{b}_{1}^{*j}}{\psi_{j}} \right), \sigma^{2} \left( \left( \sum_{j=1}^{m} \psi_{j}^{-1} \right)^{-1} \left( (\boldsymbol{X}_{1}'\boldsymbol{X}_{1})^{-1} + (\boldsymbol{X}'\boldsymbol{X})^{-1} \right) + (\boldsymbol{X}'\boldsymbol{X})^{-1} \right) \right],$$

$$\sigma^{2} | \psi_{1}, \dots, \psi_{m}, \boldsymbol{b}_{1}^{*1}, \dots, \boldsymbol{b}_{1}^{*m}, \operatorname{RSS}_{1}^{*1}, \dots, \operatorname{RSS}_{1}^{*m} \sim \operatorname{Scale-inv-} \chi^{2} \left( \nu, \tilde{\tau}_{2}^{2} \right),$$

$$\pi(\psi_{1}, \dots, \psi_{m} | \boldsymbol{b}_{1}^{*1}, \dots, \boldsymbol{b}_{1}^{*m}, \operatorname{RSS}_{1}^{*1}, \dots, \operatorname{RSS}_{1}^{*m})$$

$$\propto \left( \prod_{j=1}^{m} \psi_{j} \right)^{-\frac{n-p+r+\alpha+1}{2}} \left( \sum_{j=1}^{m} \psi_{j}^{-1} \right)^{-\frac{p}{2}} \left( 1 + \left( \sum_{j=1}^{m} \psi_{j}^{-1} \right) \right)^{-\frac{2n-2p+\alpha-1}{2}} \left\{ \nu \tilde{\tau}_{2}^{2} \right\}^{-\frac{\nu}{2}}$$

where  $\nu = m(r-2) - p + \delta + 1$  as before and

$$\nu \tilde{\tau_2}^2 = \sum_{j=1}^m \left( \boldsymbol{b}_1^{*j} - \left( \sum_{j=1}^m \psi_j^{-1} \right)^{-1} \left( \sum_{j=1}^m \frac{\boldsymbol{b}_1^{*j}}{\psi_j} \right)^{-1} \right)^{\prime} \frac{\left( (\boldsymbol{X}_1^{\prime} \boldsymbol{X}_1)^{-1} + (\boldsymbol{X}^{\prime} \boldsymbol{X})^{-1} \right)^{-1}}{\psi_j} \\ \left( \boldsymbol{b}_1^{*j} - \left( \sum_{j=1}^m \psi_j^{-1} \right)^{-1} \left( \sum_{j=1}^m \frac{\boldsymbol{b}_1^{*j}}{\psi_j} \right)^{-1} \right) + \sum_{j=1}^m \frac{\text{RSS}_1^{*j}}{\psi_j}$$

The latent variables can be sampled using Accept-Reject algorithm similarly as before, using an Inverse-Dirichlet distribution as proposal distribution.

### 6 Discussion

We have described how to perform full Bayesian analysis in the multiple regression model based on multiply imputed synthetic data. The priors used are non-informative priors. The frequentist coverage of the credible set (based on a limited simulation study that is not reported here) is sensitive to the choice of the hyperparameter and depending on the value the coverage may be significantly lower than the nominal level. This is expected due to the latent variable structure of the problem. However, more investigation is needed to understand the impact of the hyperparameters. The analysis can be extended to related cases such as where all or some of the independent variables are also sensitive or when the response is multivariate. We will consider such an investigation in the future.

## Acknowledgemnt

The authors thank Andrew Raim for reviewing the paper and providing many helpful comments and suggestions. The authors also thank Tommy Wright for his support and encouragement.

## References

- [1] A. Guin, A. Roy, and B.K. Sinha. Bayesian analysis of singly imputed partially synthetic data generated by plug-in sampling and posterior predictive sampling under the multiple linear regression model. *International Journal of Statistical Sciences*, to appear, 2021.
- [2] Martin Klein, John Zylstra, and Bimal Sinha. Finite Sample Inference for Multiply Imputed Synthetic Data under a Multiple Linear Regression Model. *Calcutta Statistical Association Bulletin*, 71(2):63–82, 2019.
- [3] M.D. Klein, T. Mathew, and B. Sinha. Noise multiplication for statistical disclosure control of extreme values in log-normal regression samples. *Journal of Privacy and Confidentiality*, 6(1), Jun. 2014.
- [4] M.D. Klein and B. Sinha. Inference for Singly Imputed Synthetic Data Based on Posterior Predictive Sampling under Multivariate Normal and Multiple Linear Regression Models. Sankhya B, 77(2):293–311, 2015.
- [5] M.D. Klein and B. Sinha. Likelihood Based Finite Sample Inference for Singly Imputed Synthetic Data Under the Multivariate Normal and Multiple Linear Regression Models. *Journal of Privacy and Confidentiality*, 7(1), Dec. 2015.
- [6] T.E. Raghunathan, J.P. Reiter, and D.B. Rubin. Multiple Imputation for Statistical Disclosure Limitation. Journal of Official Statistics, 19(1):1–16, 2003.
- [7] J.P. Reiter. Inference for Partially Synthetic, Public Use Microdata Sets. Survey Methodology, 29(2):181–188, 2003.
- [8] D.B. Rubin. Multiple Imputation for Nonresponse in Surveys. Wiley, 1987.