

A Semi-Parametric Bayesian Generalized Least Squares Estimator

Ruochen

Wu

School of
Economics Fudan
University

Melvyn

Weeks

Faculty of
Economics and
Claire College
University of
Cambridge

Abstract

In this paper we propose a semi-parametric Bayesian Generalized Least Squares estimator. In a generic GLS setting where each error is a vector, parametric GLS maintains the assumption that each error vector has the same covariance matrix. In reality however, the observations are likely to be heterogeneous regarding their distributions. To cope with such heterogeneity, a Dirichlet process prior is introduced for the covariance matrices of the errors, leading to the error distribution being a mixture of a variable number of normal distributions. Our methods let the number of normal components be data driven. Two specific cases are then presented: the semi-parametric Bayesian Seemingly Unrelated Regression (SUR) for equation systems; as well as the Random Effects Model (REM) and Correlated Random Effects Model (CREM) for panel data. A series of simulation experiments is designed to explore the performance of our methods. The results demonstrate that our methods obtain smaller posterior standard deviations than the parametric Bayesian GLS. We then apply our semi-parametric Bayesian SUR and REM/CREM methods to empirical examples.

Reference Details

CWPE 2011

Published 20 February 2020

Updated 25 September 2020

Key Words Bayesian semi-parametric, generalized least square estimator, Dirichlet process, equation system, seemingly unrelated regression, panel data, random effects model, correlated random effects model.

JEL Codes C3

Website www.econ.cam.ac.uk/cwpe

A Semi-Parametric Bayesian Generalized Least Squares Estimator

Ruo Chen Wu
School of Economics
Fudan University

Melvyn Weeks*
Faculty of Economics and Clare College,
University of Cambridge

September 25, 2020

Abstract

In this paper we propose a semi-parametric Bayesian Generalized Least Squares estimator. In a generic GLS setting where each error is a vector, parametric GLS maintains the assumption that each error vector has the same covariance matrix. In reality however, the observations are likely to be heterogeneous regarding their distributions. To cope with such heterogeneity, a Dirichlet process prior is introduced for the covariance matrices of the errors, leading to the error distribution being a mixture of a variable number of normal distributions. Our methods let the number of normal components be data driven. Two specific cases are then presented: the semi-parametric Bayesian Seemingly Unrelated Regression (SUR) for equation systems; as well as the Random Effects Model (REM) and Correlated Random Effects Model (CREM) for panel data. A series of simulation experiments is designed to explore the performance of our methods. The results demonstrate that our methods obtain smaller posterior standard deviations than the parametric Bayesian GLS. We then apply our semi-parametric Bayesian SUR and REM/CREM methods to empirical examples.

JEL Classification Code: C3

Keywords: Bayesian semi-parametric, generalized least square, Dirichlet process, equation system, seemingly unrelated regression, panel data, random effects model, correlated random effects model.

*Contact Author: Dr. M. Weeks, Faculty of Economics, University of Cambridge, Cambridge CB3 9DD, UK.
Email: mw217@econ.cam.ac.uk. We have benefited from comments provided by Oliver Linton, ...

1 Introduction

The Generalized Least Square (GLS) estimator is a family of econometric methods that have seen numerous applications in empirical economics. As pointed out by Wooldridge (2003), parametric GLS type estimators accommodate a deviation from the assumption that the errors in the model are homoskedastic and serial uncorrelated. For example, relative to the ordinary least squares regression model, GLS no longer assumes that the covariance matrix of the errors is diagonal with identical diagonal elements.

In a more general setting, each error may be a random vector, which includes some of the most popular applications of GLS. For example, the Seemingly Unrelated Regression (SUR, Zellner, 1962 and 1971) has been developed for equation systems, and widely applied to gain efficiency by exploiting the correlation between errors across equations. Similarly, in the analysis of panel data the random effects model REM recognizes that there are individual specific, time-invariant features that are unobservable and uncorrelated with the explanatory variables. A useful extension to RE is the correlated random effects (CRE) model (Chamberlain, 1980; Wooldridge, 2005; Murtazashvili and Wooldridge, 2008), which allows the individual effects to be correlated with the explanatory variables usually as a linear function of the means of the regressors.

However, the parametric GLS still maintains the assumption that the error vector of each individual has the same covariance matrix. In reality however, heterogeneity in error distributions is a major concern in empirical analysis. Such heterogeneity can be caused by observations on individuals or households reflecting variation in demographics such as the size of the household, and the level of income among others. It is a challenge for analysts who seek reliable inference with the data to capture the form of the heterogeneity in observations.

The standard Bayesian approach to GLS assumes that the error distribution is multivariate normal. Recent developments in Bayesian methods allow the use of prior information to relax this assumption. For example, the Dirichlet prior has been introduced to accommodate heterogeneity in the distributions of both errors (see Chigira and Shiba, 2015 for an example) and model parameters (Allenby et al., 1998) by mixing a fixed number of normal distributions.

However, a notable drawback of the Dirichlet prior is that the the dimension of the mixing distribution is usually unknown. Bayesian semi-parametric methods introduce more flexibility by letting the data and the prior determine the structure of heterogeneity jointly. The Dirichlet Process (\mathcal{DP}) prior¹ can be used to form a mixing of normal distributions, whose dimension need not be predetermined. In this sense, the use of \mathcal{DP} priors represents a more flexible approach to accommodating heterogeneity than the mixing of a fixed number of normal distributions with the Dirichlet prior.

In the context where heterogeneity in the distributions of errors is the major concern, as in where the focus is upon the inference, \mathcal{DP} priors are introduced for the hyper-parameters of the errors in the model² This leads to the grouping of the hyper-parameters, with those in the same group having identical values. As such, the corresponding errors of hyper-parameters in the same group have the same distribution, while the errors whose hyper-parameters are in separate groups are from different distributions.

A landmark study in this area is Conley et al. (2008), where a Bayesian semi-parametric approach to the instrumental variable problem was introduced in a two stage least square framework. Due to the endogeneity of some explanatory variables, the errors in the two stages are correlated by construction. Instead of assuming that the joint errors in the two stages have an identical bivariate normal distribution (cf. Chao and Phillips, 1998; Geweke, 1996; Kleibergen

¹See Escobar and West, 1995 and 1998 and MacEachern, 1998 for a reference of the Dirichlet Process prior.

²E.g. for an error vector with zero mean, its hyper-parameter is its covariance matrix.

and van Dijk, 1998; Rossi et al., 2005), the authors introduced a Dirichlet process prior for their hyper-parameters. This provides a semi-parametric version of two stage least squares, where the errors of the two stages jointly follow a non-parametric mixture of normal distributions.

In this paper we also focus on relaxing the identical distribution assumption on the errors, but in a different scenario from that of Conley et al. (2008). We propose a semi-parametric Bayesian GLS that incorporates the \mathcal{DP} prior. The motivation is to explore more information in the error distribution by allowing their hyper-parameters to differ across observations. The resulting distribution of the error terms will involve a mixture of normal distributions where the number of the normal components is influenced by both the prior and the data. We then introduce two specific cases of semi-parametric Bayesian GLS, namely for equation systems and panel data.

The rest of the paper is organized as follow. In Section 2 we introduce the generic form of the Dirichlet process, and demonstrate its use as a prior for semi-parametric Bayesian GLS. Then two special cases of the GLS are described. The DP-SUR method is introduced in Section 3. Sections 3.3 and 3.4 present the simulation design and results, respectively, for the DP-SUR. Two empirical examples are given in Section 3.5. Section 4 motivates and introduces our semi-parametric Bayesian GLS methods for panel data, the DP-REM, and its extension, the DP-CREM is introduced in Section 4.3. Simulation designs and results for the panel setting are in Section 4.4. The DP-REM and DP-CREM methods are then applied to two empirical examples in Section 4.5. Section 5 concludes the paper.

2 Bayesian GLS with Dirichlet Process Prior

We introduce the generic Bayesian GLS with \mathcal{DP} mixture in this section. In Section 2.1 we briefly review the literature in related areas. Then in Section 2.2 we present how the Dirichlet process prior can be used to produce a semi-parametric Bayesian GLS.

2.1 Literature Review

Bayesian attempts to incorporate heterogeneity in the hyper-parameters of the errors, and consequently in their distributions, can be traced back to Geweke (1993). He introduced a Bayesian GLS where an inverse gamma prior is introduced for variances of the errors, each of which had a normal distribution. Geweke (1993) demonstrated that such a scale mixture of normal distributions is equivalent to the errors having a t-distribution.

Although the model with the t-distributed errors is flexible, this approach depends upon the assumption that the normal distributions are mixed with inverse gamma distributed variances. As pointed out by Koop (2003), relaxing this assumption results in more flexible models, given that the errors are no longer restricted to having a t-distribution. This can be done by using a Dirichlet prior (the conjugate prior of a multinomial distribution) to mix a finite number of normal distributions. The Dirichlet mixture model has emerged as a widely applied methodology for capturing heterogeneity in both linear and non-linear models, including Allenby et al. (1998), Li and Tobias (2011) and Chigira and Shiba (2015). However, the main limitation is that it takes a fairly difficult test procedure to determine the “correct” number of mixing components.

In the wake of this limitation of the Dirichlet mixture model, it seems more reasonable to let the data and the prior determine the number of normal components jointly. This can be achieved using a Dirichlet prior of infinite dimension, which is the Dirichlet process (\mathcal{DP}) introduced by Ferguson (1973)³. \mathcal{DP} is the conjugate prior for an infinite dimension, non-

³See Teh (2011) and Gershman and Blei (2012) for reviews of the Dirichlet process.

parametric multinomial distribution. The generic form of the \mathcal{DP} can be written as

$$F \sim \mathcal{DP}(\alpha, F_0), \quad (1)$$

where $\alpha > 0$ is the concentration parameter, and F_0 is the base distribution. F is a random distribution that is discrete with probability one⁴.

The \mathcal{DP} is a non-parametric “distribution of distributions” (Escobar and West, 1995 and 1998; MacEachern, 1998), in that a draw, say F , from a \mathcal{DP} is a probability distribution itself. The Chinese restaurant process (Aldous, 1985) provides the predictive probabilities of the n^{th} realization, r_n , conditioned on the $n - 1$ existing realizations $\{r_1, r_2, \dots, r_{n-1}\}$, and realizations from F can be drawn accordingly. Due to the fact that F is discrete, the existing realizations will be assigned to groups, each with a unique value for all realizations in the same group. Denote the group id of r_i as $c_i = 1, \dots, K$, and the unique value of group c_i as $r_{c_i}^*$: if r_i is in group k , then $r_i = r_{c_i}^* = r_k^*$. The prediction probabilities of r_n is given by

$$\Pr\{r_n = r_k^* | r_1, r_2, \dots, r_{n-1}\} = \begin{cases} \frac{n_k}{n-1+\alpha} & \text{if } 1 \leq k \leq K \\ \frac{\alpha}{n-1+\alpha} & \text{if } k = K+1 \text{ (i.e. } r_n = r_{K+1}^* \sim F_0) \end{cases}, \quad (2)$$

where n_k is the number of realizations that are already in group k . Aldous (1985) showed that r_1, r_2, \dots, r_n ⁵ generated according to the Chinese restaurant process are i.i.d. draws from F , i.e.

$$\begin{aligned} F | \alpha, F_0 &\sim \mathcal{DP}(\alpha, F_0) \\ r_i | F &\stackrel{iid}{\sim} F. \end{aligned} \quad (3)$$

A model with a \mathcal{DP} prior on the *distribution* of parameters is called a \mathcal{DP} mixture model (de Carvalho et al., 2013; Wiesenfarth et al., 2014; Li et al., 2018 and Hejblum et al., 2019), and is capable of representing very general forms of heterogeneity in the distributions of the observations. The \mathcal{DP} normal mixture model, i.e. one whose mixture components are normal distributions, can be written as

$$\begin{aligned} F | \alpha, F_0 &\sim \mathcal{DP}(\alpha, F_0) \\ \theta_i | F &\stackrel{iid}{\sim} F, \\ \mathbf{y}_i | \theta_i &\sim \mathcal{N}(\theta_i), \end{aligned} \quad (4)$$

where θ_i is the set of parameters of observation \mathbf{y}_i . In the multivariate normal⁶ case, θ_i consists of the mean vector and covariance matrix, i.e. $\theta_i = (\boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i)$.

The posterior probability of θ_i having the same value as one of the existing θ_{-i} is

$$\Pr\{\theta_i = \theta_k^* | \theta_{-i}, \mathbf{y}_i, \alpha\} \propto \frac{n_k}{n-1+\alpha} \mathcal{N}(\mathbf{y}_i | \theta_k^*), \quad (5)$$

where θ_k^* and n_k denote, respectively, the unique value of group k and the number of observations already in group k . The posterior probability of θ_i taking a new value from the base distribution, i.e. $\theta_i = \theta_{new}^* \sim F_0$, is

$$\Pr\{\theta_i = \theta_{new}^* | \theta_{-i}, \mathbf{y}_i, \alpha, F_0\} \propto \frac{\alpha}{n-1+\alpha} \int \mathcal{N}(\mathbf{y}_i | \theta_{new}^*) p(\theta_{new}^* | F_0) d\theta_{new}^*, \quad (6)$$

where $p(\theta_{new}^* | F_0)$ is the probability density of the new value θ_{new}^* given F_0 . We now introduce how the \mathcal{DP} normal mixture model is applied to introduce a semi-parametric Bayesian GLS.

⁴The level of discreteness is influenced by α , the concentration parameter.

⁵The realisations r_1, r_2, \dots, r_n generated according to (2) are not independent given that the n^{th} realisation is generated conditioned on the $n - 1$ realizations before. However, these realisations are exchangeable, and therefore independent conditional on a distribution F .

⁶Note that \mathbf{y}_i is a vector.

2.2 Semi-parametric Bayesian GLS

Below we introduce the generic form of the semi-parametric GLS estimator, where a \mathcal{DP} prior is introduced on the distribution of the hyper-parameters of the errors. Consider a general linear regression

$$\mathbf{y}_i = \mathbf{X}_i\boldsymbol{\beta} + \boldsymbol{\varepsilon}_i, \quad (7)$$

where i indexes the observation, \mathbf{y}_i is a $Q \times 1$ vector of dependent variables, \mathbf{X}_i is a $Q \times K$ matrix of explanatory variables, $\boldsymbol{\beta}$ is a $K \times 1$ vector of coefficients, and $\boldsymbol{\varepsilon}_i$ is a $Q \times 1$ vector of errors. Our semi-parametric GLS estimator introduces a \mathcal{DP} prior on the distribution of the error covariance matrix which will be used to weight the observations. Given the usual assumption of zero means for the errors then $\theta_i = \boldsymbol{\Sigma}_i$. The hierarchical prior can then be written as

$$\begin{aligned} F|\alpha, F_0 &\sim \mathcal{DP}(\alpha, F_0) \\ \boldsymbol{\Sigma}_i|F &\stackrel{iid}{\sim} F. \end{aligned} \quad (8)$$

Due to the discreteness of F under the \mathcal{DP} prior, the value of some covariance matrices $\boldsymbol{\Sigma}_i$ will be the same, thus putting $\boldsymbol{\Sigma}_i$ into groups denoted by c_i . This ‘‘grouping’’ characteristic can help to reveal the structure of the unobserved heterogeneity in the data.

The GLS estimator weights the observations according to their covariance matrix. The likelihood of $\boldsymbol{\beta}$ is then

$$p(\mathbf{y}_i|\boldsymbol{\beta}, \boldsymbol{\Sigma}_{c_i}^*) = \frac{1}{(2\pi)^{Q/2}} |\boldsymbol{\Sigma}_{c_i}^*|^{-\frac{1}{2}} \exp \left[-\frac{1}{2} (\mathbf{y}_i - \mathbf{X}_i\boldsymbol{\beta})' \boldsymbol{\Sigma}_{c_i}^{*-1} (\mathbf{y}_i - \mathbf{X}_i\boldsymbol{\beta}) \right]. \quad (9)$$

where $\boldsymbol{\Sigma}_{c_i}^*$ is the unique covariance matrix of group c_i . Given the choice of prior for $\boldsymbol{\beta}$, one could generate draws from the posterior of the parameters with MCMC methods.

For the conjugate normal prior for $\boldsymbol{\beta}$ is specified, i.e.

$$\boldsymbol{\beta} \sim \mathcal{N}(\mathbf{b}_0, \mathbf{V}_0), \quad (10)$$

where \mathbf{b}_0 and \mathbf{V}_0 denote, respectively, the prior mean and covariance matrix of $\boldsymbol{\beta}$. The posterior of $\boldsymbol{\beta}$ may then be written as

$$\boldsymbol{\beta}|\mathbf{y}, \boldsymbol{\Sigma}_{c_i}^* \sim \mathcal{N}(\mathbf{b}, \mathbf{V}), \quad (11)$$

where

$$\mathbf{V} = \left(\mathbf{V}_0^{-1} + \sum_{i=1}^N \mathbf{X}_i' \boldsymbol{\Sigma}_{c_i}^{*-1} \mathbf{X}_i \right)^{-1}, \quad (12)$$

and

$$\mathbf{b} = \mathbf{V} \left(\mathbf{V}_0^{-1} \mathbf{b}_0 + \sum_{i=1}^N \mathbf{X}_i' \boldsymbol{\Sigma}_{c_i}^{*-1} \mathbf{y}_i \right). \quad (13)$$

Note that (11), (12) and (13) have the same form as the posterior of the parametric Bayesian GLS estimator assuming *i.i.d.* normal errors. In the case of the semi-parametric Bayesian GLS estimator, the errors are associated with different hyper-parameters, such that each observation i is weighted by $\boldsymbol{\Sigma}_{c_i}^*$. DP-GLS is generic, without making any assumptions on the form of the covariance matrix other than all $\boldsymbol{\Sigma}_{c_i}^*$ being positive definite and symmetric. We now proceed to explain how the semi-parametric Bayesian GLS estimators work in two specific contexts, which are equation systems in Section 3 and panel data models in Section 4.

3 Semi-parametric Seemingly Unrelated Regression

Below we introduce the SUR equation system and demonstrate how the \mathcal{DP} prior is incorporated. Without loss of generality we consider a system of two equations

$$\begin{aligned} y_{1i} &= \beta_{10} + x_{11,i}\beta_{11} + x_{12,i}\beta_{12} + \varepsilon_{1i} \\ y_{2i} &= \beta_{20} + x_{21,i}\beta_{21} + x_{22,i}\beta_{22} + x_{23,i}\beta_{23} + \varepsilon_{2i}, \end{aligned} \quad (14)$$

where y_{mi} denotes observation i for equation m ($m = 1, 2$) and $x_{mk,i}$ ($k = 1, 2, 3$) are the explanatory variables. β_{ml} ($l = 0, 1, 2, 3$) denote the coefficients, and ε_{mi} are the errors. In this context, the dimension of each error vector, i.e. Q in Section 2.2 will be the number of equations in the system.

The model can be written in matrix form

$$\begin{aligned} \mathbf{y}_1 &= \mathbf{X}_1\boldsymbol{\beta}_1 + \boldsymbol{\varepsilon}_1 \\ \mathbf{y}_2 &= \mathbf{X}_2\boldsymbol{\beta}_2 + \boldsymbol{\varepsilon}_2, \end{aligned} \quad (15)$$

where $\mathbf{y}_m = \{y_{mi}\}$, $\boldsymbol{\varepsilon}_m = \{\varepsilon_{mi}\}$ are $N \times 1$ vectors. $\mathbf{X}_1 = [\boldsymbol{\iota}, \mathbf{x}_{11}, \mathbf{x}_{12}]$ and $\mathbf{X}_2 = [\boldsymbol{\iota}, \mathbf{x}_{21}, \mathbf{x}_{22}, \mathbf{x}_{23}]$ are $N \times 3$ and $N \times 4$ matrices, respectively, where $\boldsymbol{\iota}$ is an $N \times 1$ vector of ones. $\boldsymbol{\beta}_1 = \{\beta_{1l}\}$ and $\boldsymbol{\beta}_2 = \{\beta_{2l}\}$ are 3×1 and 4×1 vectors, respectively.

In the presence of correlated errors there exists an efficiency gain by utilising a system estimator. The Seemingly Unrelated Regression (SUR, Zellner, 1962) was introduced for this task. Instead of $\varepsilon_{1i} \stackrel{iid}{\sim} \mathcal{N}(0, \sigma_1^2)$ and $\varepsilon_{2i} \stackrel{iid}{\sim} \mathcal{N}(0, \sigma_2^2)$, as in the OLS case, the errors $\boldsymbol{\varepsilon}_i$ are now identically multivariate normally distributed, i.e. $\boldsymbol{\varepsilon}_i = (\varepsilon_{1i} \ \varepsilon_{2i})' \stackrel{iid}{\sim} \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma})$. The covariance matrix of $\boldsymbol{\varepsilon}$ is then

$$\boldsymbol{\Omega} = \boldsymbol{\Sigma} \otimes \mathbf{I} = \begin{bmatrix} \sigma_{11}\mathbf{I}_N & \sigma_{12}\mathbf{I}_N \\ \sigma_{21}\mathbf{I}_N & \sigma_{22}\mathbf{I}_N \end{bmatrix}, \text{ s.t. } \sigma_{12} = \sigma_{21}, \quad (16)$$

where " \otimes " stands for the Kronecker product.

One could transform the observations with this covariance matrix, so that the errors follow the standard normal distribution $\mathcal{N}(0, 1)$, with the likelihood, prior and posterior of the parameters defined similarly as in equations (9) to (13)⁷.

3.1 DP Prior for SUR

Although the SUR model accounts for the cross-equation correlation of errors, as Wooldridge (2003) has noted, the errors are assumed to be identically distributed. Moreover, unlike the classical GLS estimator, this distribution is usually assumed to be normal. In this section we propose a new DP-SUR method that makes no a priori assumptions on the family of distribution of the errors. If we allow each observation i have its own covariance matrix, flexibility of the error distribution will lead to identification problems if we only have a cross sectional data. Assigning the observations into groups represents a compromise. Given (15), the covariance matrix of the error for observation i is given by

$$\boldsymbol{\Sigma}_i = \begin{bmatrix} \sigma_{11,c_i}^* & \sigma_{12,c_i}^* \\ \sigma_{21,c_i}^* & \sigma_{22,c_i}^* \end{bmatrix}, \text{ s.t. } \sigma_{12,c_i}^* = \sigma_{21,c_i}^*, \quad (17)$$

where c_i denotes the group id of observation i , and the superscript $*$ denotes the group-specific hyper-parameter. For $c_i = c_j$, $i, j \in \{1, 2, \dots, N\}$ observations i and j share the same group ID

⁷However, the covariance matrix $\boldsymbol{\Omega}$ has a specific form as in the SUR in (16), instead of the general, positive definite symmetric form of a covariance matrix.

and hyper-parameters, such that $\sigma_{pq,c_i}^* = \sigma_{pq,c_j}^*$, where $p, q \in \{1, 2\}$ index the equations in the system.

Assuming that the number of groups are known, the Dirichlet prior may be used to perform the mixing. A less restrictive approach utilises a non-parametric approach by introducing a \mathcal{DP} prior for the distribution of Σ_i as in (8). A natural choice of the base distribution F_0 is the conjugate prior for the covariance matrix of a multivariate normal distribution, the inverse Wishart distribution, i.e.

$$F_0 \equiv \mathcal{IW}(\nu, \mathbf{W}) \quad (18)$$

where ν, \mathbf{W} are hyper-parameters of the inverse Wishart distribution. Given (18) the posterior distribution of each Σ_i is also inverse Wishart, which is easy to draw from using the Gibbs sampler. The main difference from the parametric Bayesian SUR is that the covariance matrix of each observation is now given in (17), which allows each group of observations to have its own unique values for the parameters.

3.2 MCMC Algorithm

A Gibbs sampler (see Geweke, 1996) is available for the SUR model. The Gibbs sampler draws two sets of parameters from their posteriors: the covariance matrix of the errors Σ and the regression parameters β , namely

$$\Sigma | \mathbf{y}, \mathbf{X}, \beta \quad (19)$$

$$\beta | \mathbf{y}, \mathbf{X}, \Sigma. \quad (20)$$

When introducing the hierarchical structure which includes the \mathcal{DP} prior, a number of extra parameters are included in the MCMC algorithm. These are the covariance matrices of the errors, $\Theta = \{\Sigma_i\}$, and α , the concentration parameter of the \mathcal{DP} prior. The Gibbs sampler now consists of

$$\Theta | \mathbf{y}, \mathbf{X}, \beta, \alpha \quad (21)$$

$$\beta | \mathbf{y}, \mathbf{X}, \Theta, \alpha \quad (22)$$

$$\alpha | \mathbf{y}, \mathbf{X}, \beta, \Theta. \quad (23)$$

The major difference between the two Gibbs sampler lies in (19) and (21). In (19) the errors have the same covariance matrix Σ . In contrast, there will be $K \leq N$ unique values in Θ in equation (21) due to the discreteness of F under the \mathcal{DP} prior; observations with the same value of Σ_i are assigned to the same group. With the last draw of β , the residuals can be obtained, which are used as the data to take a draw for Θ .

In making draws of the concentration parameter α using (23), we adopt the \mathcal{DP} prior introduced by Conley et al. (2008), namely

$$p(\alpha) \propto \left(1 - \frac{\alpha - \alpha_{min}}{\alpha_{max} - \alpha_{min}} \right)^\tau, \quad (24)$$

where α_{min} and α_{max} are the pre-set lower and upper bound of α . Larger α lead to more groups being generated on average, i.e. the \mathcal{DP} being less discrete. According to the distribution of the number of groups K conditioned on α in Antoniak (1974), we could determine α_{min} and α_{max} by setting the mode of number of groups to K_{min} and K_{max} . In this paper we let K_{min} be 1 and K_{max} be 5% of the number of observations. Following the suggestion of Conley et al. (2008), we set τ to 0.8. The hyper-parameters α_{max} has been adjusted according to K_{max} being 10% and 50% of the sample size. In our experiments the results are insensitive to these changes in the hyper-parameters in the prior of the concentration parameter α .

3.3 A Simulation Experiment

In this section we conduct a simulation experiments designed to compare our method to the Bayesian SUR described in Section 3. As the main focus of this paper is the potential efficiency gains over GLS type estimators, we evaluate the performance of the DP-SUR and normal Bayesian SUR focusing upon the posterior standard deviations of the parameters estimated with the two methods. All simulation experiments are based upon the two equation system in (14).

The experiments are designed to highlight the performance of the estimators along the following dimensions:

- (i) heterogeneity in the errors;
- (ii) the tail of the error distribution;
- (iii) sample size.

For (i) we check the performance of our DP-SUR approach against a model where the errors are distributed *i.i.d.* multivariate normal. In the heterogeneous case, the most direct way is to generate the errors from a mixture of multivariate normal distributions⁸. However, we use multivariate t-distributions (Andrews and Mallows, 1974) to exploit the scale mixtures of normal distributions with inverse Wishart covariance matrices.

To accommodate (ii), we vary the degrees of freedom (df) of the multivariate t-distribution. Smaller degrees of freedom leads to heavier tails, which indicates that a larger proportion of observations follow normal distributions that are “flatter”, i.e. less concentrated around the mean.

To determine the robustness of our method, we include a set of simulations where the errors follow a log-normal distribution. The log-normal distribution has seen a wide range of applications in empirical studies. For example, with perhaps the exception of the top 1-3 percent of the population, income has been shown to follow a log-normal distribution (Clementi and Gallegati, 2005). In addition, extreme realizations are more likely to be generated from the multivariate log-normal distribution, as it is fat-tailed.

Using (15), the explanatory variables are drawn from normal distributions with parameters

$$x_{11,i} \stackrel{iid}{\sim} \mathcal{N}(1, 1), \quad x_{12,i} \stackrel{iid}{\sim} \mathcal{N}(3, 1),$$

and

$$x_{21,i} \stackrel{iid}{\sim} \mathcal{N}(-2, 1), \quad x_{22,i} \stackrel{iid}{\sim} \mathcal{N}(4, 1), \quad x_{23,i} \stackrel{iid}{\sim} \mathcal{N}(-1, 1). \quad (25)$$

We set $\beta_1 = (1.0, -0.5, 1.6)'$ and $\beta_2 = (1.5, -1.2, -0.7, 2)'$. We generate errors from the multivariate normal distribution $\mathcal{N}(\mathbf{0}, \Sigma)$, where

$$\Sigma = \begin{bmatrix} 1 & 0.5 \\ 0.5 & 1 \end{bmatrix}. \quad (26)$$

Without loss of generality, we let the variances be identical, and fix the correlation between the errors in the two equations at 0.5. We set three samples sizes for the simulation experiment: 100, 250, and 500.

When generating errors from a multivariate t-distribution, we set the location parameter to $\mu = \mathbf{0}$ and shape parameter to Σ . As noted, the parameter that controls the tail behaviour of

⁸Simulating data from a mixture of multivariate normal distributions can be problematic given the influence of the following: the number of components, the covariance matrices (we fix the mean at zero) of the normal components and the weights assigned to each component.

multivariate t-distributions is the df. We set the df to 2⁹, 3 and 4. For df=2 the tails of the corresponding multivariate t-distribution are much heavier than that of the multivariate normal with the same location and shape parameters $\boldsymbol{\mu}$ and $\boldsymbol{\Sigma}$. When the df is 4, the tails of the multivariate t-distribution are only slightly heavier than the multivariate normal. Our DP-SUR should demonstrate relative efficiency in all three situations, as the multivariate t-distributions has heavier tails than the multivariate normal. Gains in efficiency will be decreasing in df, given that the tails are less heavy.

The errors in the multivariate log-normal scenario are constructed by first drawing from multivariate normal distributions, and then taking the natural exponent of these draws. Because the log-normal distribution has a positive mean, it is necessary to demean, so that the errors will have zero means. Our DP-SUR is expected to have efficiency gains in the log-normal case given it is asymmetric and heavy tailed.

3.4 DP-SUR Simulation Results

Below we present the simulation results¹⁰. We present the posterior standard deviations (s.d.) estimated with both our DP-SUR method and the Bayesian SUR assuming multivariate normal errors, along with the percentage differences between them¹¹. As the behaviour of the posteriors of all the parameters are uniformly similar, for the sake of clearness, we present only the results regarding β_{11} and β_{21} ¹².

Multivariate t-distributed Errors

Table 1 presents the posterior s.d.s and the percentage difference¹³ between the s.d. estimated with the DP-SUR and the normal SUR, both averaged over the samples. We observe that the DP-SUR gives smaller posterior s.d when df is 2, 3 and 4. The percentage differences when df is 2 are above 40%, above 20% when df is 3, and around 15% when df is 4 as shown in the upper three panels of the table. Efficiency gains increase with sample size as more extreme values of the errors are realised. The parametric SUR assumes that all the realizations have the same $\boldsymbol{\Sigma}$, where the extreme ones will expand this $\boldsymbol{\Sigma}$ shared by all realizations. In contrast, these extreme realizations will be assigned to distributions with larger $\boldsymbol{\Sigma}_i$'s by the DP-SUR, while the rest will be treated as realizations from normal distributions with smaller $\boldsymbol{\Sigma}_i$'s. By accommodating a higher degree of heterogeneity in the error distributions, potentially more efficiency gains could be achieved by the DP-SUR.

Our results are consistent with expectations. The efficiency gains of the semi-parametric DP-SUR are the largest when the df is 2 (with the heaviest tails). Efficiency gains fall with the df increasing, given less heavy tails of the distribution of the errors. In fact, the lowest panel in Table 2 where the df is infinity, we observe that the posterior s.d. estimated with the two methods are very close. The s.d. for DP-SUR is slightly larger than their SUR counterparts. This is not surprising since when the distribution of the errors is multivariate normal, the parametric method is more parsimonious, using the correct structure for the covariance matrix of the errors. In the multivariate normal case, among the three sample sizes, the differences between the s.d. are the largest¹⁴ when the sample size is 100. This is expected as the information “wasted” by

⁹We do not use df 1, as the t distribution does not even have a mean in this case.

¹⁰We carry out 100 simulations for each sample size, which proved sufficient to achieve stable results even with the smallest sample size.

¹¹The tables containing the posterior means can be found in the Appendices. For the tables of posterior means, there are 6 columns presenting the means estimated by the two methods for the 3 correlations.

¹²The full results are included in Appendices.

¹³ $\Delta\% = (\text{s.d.}_{\text{SUR}} - \text{s.d.}_{\text{DP}}) / \text{s.d.}_{\text{SUR}} \times 100\%$.

¹⁴Nevertheless, the differences are still small in magnitude, less than 2.5% for all coefficients.

the less parsimonious DP-SUR in this case has a larger impact on efficiency when the sample size is small.

Table 1: Posterior S.D., multivariate t errors

df = 2									
Sample size	100			250			500		
Parameters	DP	SUR	$\Delta\%$	DP	SUR	$\Delta\%$	DP	SUR	$\Delta\%$
β_{11}	0.1149	0.2155	41.57%	0.0732	0.1596	49.56%	0.0458	0.1126	54.22%
β_{21}	0.1107	0.2080	41.49%	0.0649	0.1359	48.48%	0.0508	0.1198	52.78%
df = 3									
Sample size	100			250			500		
Parameters	DP	SUR	$\Delta\%$	DP	SUR	$\Delta\%$	DP	SUR	$\Delta\%$
β_{11}	0.1114	0.1452	21.20%	0.0708	0.0973	26.05%	0.0446	0.0634	28.53%
β_{21}	0.1079	0.1397	20.63%	0.0630	0.0866	25.91%	0.0488	0.0694	28.70%
df = 4									
Sample size	100			250			500		
Parameters	DP	SUR	$\Delta\%$	DP	SUR	$\Delta\%$	DP	SUR	$\Delta\%$
β_{11}	0.1079	0.1248	13.52%	0.0696	0.0826	15.22%	0.0440	0.0529	16.58%
β_{21}	0.1054	0.1200	12.15%	0.0618	0.0733	15.08%	0.0473	0.0582	18.41%
df = ∞									
Sample size	100			250			500		
Parameters	DP	SUR	$\Delta\%$	DP	SUR	$\Delta\%$	DP	SUR	$\Delta\%$
β_{11}	0.0854	0.0837	-2.09%	0.0569	0.0565	-0.85%	0.0372	0.0371	-0.56%
β_{21}	0.0879	0.0860	-2.25%	0.0575	0.0572	-0.58%	0.0388	0.0384	-1.14%

Multivariate Log-normal Errors

The posterior S.D.s are presented in Table 2. We observe that the DP-SUR posterior S.D. are more than 55% smaller than those calculated using the Bayesian SUR assuming *i.i.d.* normal errors. The efficiency gains increase with sample size, which reach more than 65% in the case of 500 observations. As with the case of t distributed errors, this is due to the fact that more extreme realizations of errors are present in larger samples, leading to more efficiency gains by grouping them.

Table 2: Posterior S.D., multivariate log-normal errors

Log-normal									
Sample size	100			250			500		
Parameters	DP	SUR	$\Delta\%$	DP	SUR	$\Delta\%$	DP	SUR	$\Delta\%$
β_{11}	0.0720	0.1831	57.64%	0.0415	0.1266	65.06%	0.0272	0.0834	66.61%
β_{21}	0.0800	0.2119	58.78%	0.0439	0.1277	64.38%	0.0270	0.0826	66.39%

3.5 DP-SUR Empirical Examples

Below we apply our DP-SUR method to an economic model of the demand for factors of production with a generalized Leontief cost function (Diewert, 1971), an equation system with the

number of equations as that of factors.¹⁵ To make our empirical demonstration as general as possible, we do not impose symmetry or homogeneity restrictions.

The dataset, taken from Malikov et al. (2016), contains 2397 observations on 285 large U.S. banks between 2001 and 2010. The data includes quantities and prices of the inputs, i.e. labour, physical assets and borrowed funding, and the quantity of output, which is the loans made by a bank. Given the relatively large sample size, it is possible for us to explore the performance of the DP-SUR with different sample sizes.

The demand for factors equation system may be written as

$$a_L = \frac{L}{Y} = \beta_{LL} + \beta_{LA} \frac{P_A}{P_L} + \beta_{LF} \frac{P_F}{P_L} + \beta_{LT} T + \varepsilon_L \quad (27)$$

$$a_A = \frac{A}{Y} = \beta_{AA} + \beta_{AL} \frac{P_L}{P_A} + \beta_{AF} \frac{P_F}{P_A} + \beta_{AT} T + \varepsilon_A \quad (28)$$

$$a_F = \frac{F}{Y} = \beta_{FF} + \beta_{FL} \frac{P_L}{P_F} + \beta_{FA} \frac{P_A}{P_F} + \beta_{FT} T + \varepsilon_F, \quad (29)$$

where L , A and F denote the quantity of labour, physical assets and borrowed funds, respectively; T denotes the trend variable; Y denotes output, and P_k is the price of factor k , with $k \in \{L, A, F\}$. For the errors we assume $(\varepsilon_{Li}, \varepsilon_{Ai}, \varepsilon_{Fi}) \sim \mathcal{N}(\mathbf{0}, \mathbf{\Sigma}_i)$ where $\mathbf{\Sigma}_i$ is the covariance matrix of observation i . We allow the errors to be correlated across the three equations in the system, i.e. $\text{cov}(\varepsilon_{ki}, \varepsilon_{si}) \neq 0$, with $k, s \in \{L, A, F\}$ indexing equations.

Given that the main objects of interest are the price elasticities, we report the posterior means and s.d. of the price elasticities of the three factors. With the generalized Leontief cost function, the cross price elasticities of the factors are given by

$$e_{ks} = \frac{1}{2} \frac{\beta_{ks} (P_k/P_s)^{-1/2}}{a_k}, \quad \forall k \neq s, \quad (30)$$

The own price elasticities are

$$e_{kk} = -\frac{1}{2} \frac{\sum_{s \neq k} \beta_{ks} (P_k/P_s)^{-1/2}}{a_k}. \quad (31)$$

Table 3 contains the posterior means of the price elasticities of the demand for factors. One can see that with both the 800 observation sub-sample and the full sample, the posterior means of all the elasticities are relatively small indicating that the demand for factors (labour, physical assets and borrowed fundings) of U.S. banks are relatively price inelastic.

Note that the own price elasticities of labour and physical assets are negative in both samples. In contrast, the own price elasticity of borrowed fundings is positive, although we note that the absolute values are extremely small¹⁶ compared to those of the labour and physical assets. This shows that the demand for borrowed fundings is inelastic in the production of the U.S. banking industry. One potential reason is that borrowed funds are usually used in ways such as to meet the required reserve ratio set by the Fed. Such a feature makes borrowed fundings rather inelastic with respect to their price.

There are some differences between the posterior means estimated with the DP-SUR and the Bayesian SUR assuming normal errors. Such differences are not observed in the simulation studies. However, it should be noted that in the simulations, the regression equation was correctly specified, which is not guaranteed with the empirical data. Such differences in the posterior

¹⁵Note that the SUR and OLS estimators are exactly the same when all the equations in the system share the same explanatory variables.

¹⁶The posterior s.d. are also relatively large for this elasticity as shown in Table 4.

means with empirical datasets have also been observed in the literature on semi-parametric mixture with a \mathcal{DP} prior, including Conley et al. (2008).

Table 3: Elasticities, U.S. banking industry: posterior means

800-observation sub-sample						
Input	Labour		Assets		Fundings	
Parameters	DP	SUR	DP	SUR	DP	SUR
Wage	-0.189	-0.422	0.375	0.254	-0.016	-0.007
Asset Price	-0.009	0.131	-0.690	-0.585	0.012	0.002
Funding Price	0.198	0.291	0.315	0.331	0.004	0.004
Full sample, 2397 observations						
Input	Labour		Assets		Fundings	
Parameters	DP	SUR	DP	SUR	DP	SUR
Wage	-0.149	-0.201	0.406	0.385	0.004	-0.001
Asset Price	-0.094	-0.068	-0.686	-0.668	-0.016	-0.005
Funding Price	0.243	0.270	0.280	0.283	0.012	0.006

Table 4 presents the posterior S.D. of the price elasticities estimated with the two samples. We observe that the DP-SUR achieves smaller posterior S.D. for all the price elasticities than the Bayesian SUR assuming normality. This is not unexpected, as the elasticities are functions of the regression parameters in the equation system, which are estimated with smaller posterior S.D. with the semi-parametric DP-SUR than the parametric Bayesian SUR. The greatest percentage difference ($\Delta\%$) with the 800-observation sub-sample takes place with the cross price elasticity of the demand for labour with respect to the price of physical assets, for which the DP-SUR posterior S.D. is 38.27% smaller than the SUR counterpart. With the full sample, the largest percentage difference is observed with the cross price elasticity of the demand for borrowed fundings with respect to the price of physical assets, which reached 39.15%.

Table 4: Elasticities, U.S. banking industry: posterior S.D.

800-observation sub-sample									
Input	Labour			Assets			Fundings		
Price	DP	SUR	$\Delta\%$	DP	SUR	$\Delta\%$	DP	SUR	$\Delta\%$
Wage	0.0341	0.0541	36.99%	0.0438	0.0484	9.60%	0.0276	0.0305	9.52%
Asset Price	0.0228	0.0369	38.27%	0.0304	0.0359	15.22%	0.0160	0.0230	30.59%
Funding Price	0.0236	0.0381	38.02%	0.0267	0.0317	15.82%	0.0148	0.0160	7.16%
Full sample, 2397 observations									
Input	Labour			Assets			Fundings		
Price	DP	SUR	$\Delta\%$	DP	SUR	$\Delta\%$	DP	SUR	$\Delta\%$
Wage	0.0189	0.0222	14.97%	0.0198	0.0210	6.01%	0.0102	0.0147	30.74%
Asset Price	0.0106	0.0131	19.38%	0.0143	0.0151	5.21%	0.0067	0.0110	39.15%
Funding Price	0.0154	0.0178	13.68%	0.0157	0.0161	2.48%	0.0081	0.0098	17.43%

In Figure 1 we present the histograms for the posterior distribution for the cross price elasticity of funding with respect to the price of physical assets for the DP-SUR and the parametric Bayesian SUR. Using the smaller 800-observation sub-sample, the posteriors for both estimators include 0. However, with the full sample the DP-SUR gives a posterior distribution that has a 95% credible interval (from -0.027 to -0.002) that excludes 0, shown by the two red vertical lines in the left panel. In contrast, the right panel shows that the parametric Bayesian SUR gives a 95% credible interval (from -0.027 to 0.016) that still includes 0 even with the full sample. From

Figure 1 we also note that the posterior distribution with the DP-SUR is strongly right skewed, which can result in the parametric Bayesian SUR having larger posterior standard deviation.

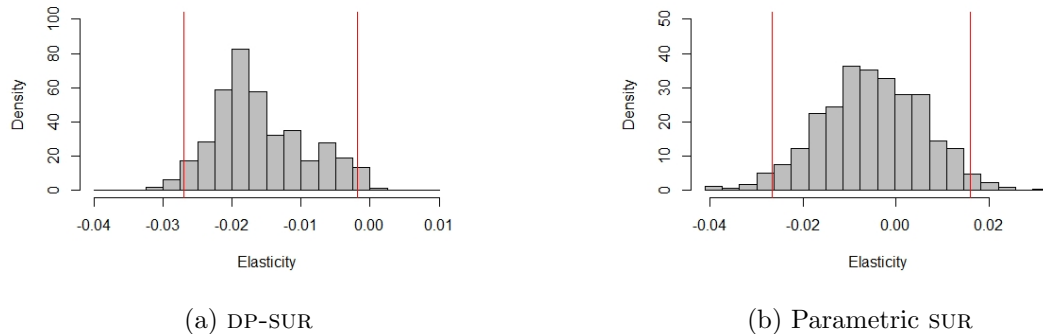


Figure 1: Histograms of the elasticity of fundings w.r.t. asset price

4 Semi-parametric Approach to Random Effects Model

In addition to equation systems, the random effects model (REM) for panel data is another scenario where the GLS has seen numerous applications. In a panel with N cross sections and time series of dimension T , the error of each individual¹⁷ is a $T \times 1$ ¹⁸ vector. We will relax the assumption of parametric Bayesian GLS for the REM (Koop, 2003) that the error vectors for all individuals have the same distribution. In this section we propose a semi-parametric Bayesian approach by introducing \mathcal{DP} priors on the variances of the random effects and the errors. We follow the same approach as in the DP-SUR method in terms of applying the \mathcal{DP} prior on the hyper-parameters.

Consider the following panel data model

$$y_{it} = \beta_1 x_{1it} + \dots + \beta_K x_{Kit} + u_i + \eta_{it} = \beta_1 x_{1it} + \dots + \beta_K x_{Kit} + \varepsilon_{it}, \quad (32)$$

where i and t index the cross section and time series dimensions of the data, respectively, y_{it} is the dependent variable, x_{kit} denote the explanatory variables, and the β_k , $k = 1, \dots, K$ are the coefficients. u_i is the time-invariant unobservable of individual i , and η_{it} the error term.

In Bayesian methods the difference between the fixed and random effects lies in the choice of prior for the individual effects u_i . Fixed effects Bayesian methods assume a non-hierarchical prior for u_i , while for the random effects a hierarchical prior is assumed. The prior for u_i may be written as

$$u_i | d^2 \stackrel{iid}{\sim} \mathcal{N}(0, d^2), \quad (33)$$

where d^2 is the variance¹⁹ of u_i . Assuming $\eta_{it} \stackrel{iid}{\sim} \mathcal{N}(0, \sigma^2)$, the posterior distribution of u_i is given by

$$u_i | \mathbf{y}_i, \boldsymbol{\beta}, d^2, \sigma^2 \sim \mathcal{N}(\mu_i, s^2), \quad (34)$$

¹⁷We use the term ‘‘individual’’ to denote the cross section unit here. In practice it can be households, firms, countries or actual individuals.

¹⁸That is, the Q in the generic semi-parametric GLS in Section 2.2 is T in this context. We use T here following panel data protocols.

¹⁹Note that the variances of u_i and η_{it} is often assumed to be random, and have their own priors. For the moment we leave them fixed for the sake of simplicity.

where $\mu_i = s^2 \sigma^{-2} \iota_T' (\mathbf{y}_i - \mathbf{X}_i \boldsymbol{\beta})$, $s^2 = (d^{-2} + \sigma^{-2} \iota_T' \iota_T)^{-1}$, with ι_T denoting a $T \times 1$ unit vector. $\mathbf{X}_i = [x_{1it}, \dots, x_{Kit}]$ is a $T \times K$ matrix of explanatory variables, and $\mathbf{y}_i = [y_{i1}, y_{i2}, \dots, y_{iT}]'$ is a $T \times 1$ vector.

The likelihood of $\boldsymbol{\beta}$ marginalized over u_i in the Bayesian REM may be written as

$$p(\mathbf{y}_i | \boldsymbol{\beta}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{T/2}} |\boldsymbol{\Sigma}|^{-\frac{1}{2}} \exp \left[-\frac{1}{2} (\mathbf{y}_i - \mathbf{X}_i \boldsymbol{\beta})' \boldsymbol{\Sigma} (\mathbf{y}_i - \mathbf{X}_i \boldsymbol{\beta}) \right], \quad (35)$$

where $\boldsymbol{\Sigma}$ is the covariance matrix of the $T \times 1$ vector $\boldsymbol{\varepsilon}_i = [\varepsilon_{i1}, \varepsilon_{i2}, \dots, \varepsilon_{iT}]'$. Assuming that $\mathbb{E}[\eta_{it} u_j | \mathbf{X}] = 0$, $\forall i, j, t$ (Greene, 2012), the covariance matrix of the composite error $\boldsymbol{\varepsilon}_i$ is

$$\text{Cov}(\boldsymbol{\varepsilon}_i) = \boldsymbol{\Sigma} = \sigma^2 \mathbf{I}_{T \times T} + s^2 \iota_T \iota_T' = \begin{bmatrix} \sigma^2 + s^2 & s^2 & \dots & s^2 \\ s^2 & \sigma^2 + s^2 & \dots & s^2 \\ \vdots & \vdots & \ddots & \vdots \\ s^2 & s^2 & \dots & \sigma^2 + s^2 \end{bmatrix}, \quad (36)$$

where σ^2 is the variance of η_{it} , and s^2 is the variance of u_i .

4.1 DP Prior for REM

Before we proceed to our DP-REM method, we review the work of Kleinman and Ibrahim (1998) and Kyung et al. (2010), who use the Dirichlet process prior for a different purpose.²⁰ Consider the model

$$y_{it} = \mathbf{X}_{it} \boldsymbol{\beta}_i + \zeta_{it}, \quad (37)$$

where $\boldsymbol{\beta}_i$ is the vector of parameters. In the literature in this area, the focus has been on the heterogeneity in the parameters (i.e. $\boldsymbol{\beta}_i$) across the individuals. For this purpose, the \mathcal{DP} prior is put on the parameters $\boldsymbol{\beta}_i$ themselves, i.e.

$$\begin{aligned} F &\sim \mathcal{DP}(\alpha, \mathcal{N}(\boldsymbol{\mu}_\beta, \boldsymbol{\Sigma}_\beta)) \\ \boldsymbol{\beta}_i | F &\stackrel{iid}{\sim} F. \end{aligned} \quad (38)$$

Given that $\boldsymbol{\beta}_i$ has a discrete DP posterior, $\boldsymbol{\beta}_i$ are grouped, with those in the same group having the same value.

Heterogeneity in the parameters *per se* is not the principle focus of this paper. Rather, this paper aims at providing more efficient inferences by exploiting the information in the distribution of unobservables. In a REM setting, unobservables are composed of individual specific unobserved effects u_i and idiosyncratic errors η_{it} . Therefore, we focus on the heterogeneity in hyper-parameters of these unobservables instead of the model parameters. In this sense, our method is in the same spirit as the literature pioneered by Conley et al. (2008).

We relax the identically distributed assumption for η_{it} and u_i , by introducing \mathcal{DP} priors on the variances. This will have the effect of grouping errors over both the cross section dimension i and the time series dimension t , with those in the same group sharing the same hyper-parameter.

The \mathcal{DP} prior for the variance of the idiosyncratic error η_{it} is

$$\begin{aligned} G &\sim \mathcal{DP}(\alpha_\eta, G_0) \\ \sigma_{it}^2 | G &\sim G, \end{aligned} \quad (39)$$

where α_η and G_0 denote the concentration parameter and base distribution of the \mathcal{DP} prior, respectively. The grouping of these variances will take place without imposing any restrictions.

²⁰The random effects model in these researches, mostly in statistics, means different from that in econometrics, as theirs is in fact random coefficients model.

For example, for $c_{it} \neq c_{is}$ ²¹ ($t \neq s$), then $\sigma_{c_{it}}^{*2}$ and $\sigma_{c_{is}}^{*2}$ are allocated to different groups and η_{it} and η_{is} have different distributions.

The \mathcal{DP} prior for the variance of the individual effects u_i in our DP-REM can be written using the following hierarchical structure

$$\begin{aligned} F &\sim \mathcal{DP}(\alpha_u, F_0) \\ d_i^2 | F &\sim F. \end{aligned} \quad (40)$$

d_i^2 is the prior variance of the random effects u_i , α_u is the concentration parameter, and F_0 the base distribution of the \mathcal{DP} prior. The use of an independent \mathcal{DP} prior on the hyper-parameters of individual effects u_i generates groupings over the N individuals such that that u_i that belong to the same group are generated from a distribution with the same hyper-parameter. This then relaxes the REM assumption that the individual effects are identically distributed.

Although the u_i are no longer identically distributed, for each particular u_i a conjugate normal prior²² can be introduced. The posterior of each u_i is then a normal distribution, the means and variances of which are different across the cross section i , i.e.

$$u_i | \mathbf{y}, \boldsymbol{\beta}, d_{c_i}^{*2}, \sigma_{c_{it}}^{*2} \sim \mathcal{N}(\mu_i, s_i^2), \quad (41)$$

where

$$\mu_i = s_i^2 \iota_T' \boldsymbol{\Sigma}_{\boldsymbol{\eta}_i}^{-1} (\mathbf{y}_i - \mathbf{X}_i \boldsymbol{\beta}) \quad (42)$$

is the posterior mean of u_i . The posterior variance is given by

$$s_i^2 = (d_{c_i}^{*-2} + \iota_T' \boldsymbol{\Sigma}_{\boldsymbol{\eta}_i}^{-1} \iota_T)^{-1}. \quad (43)$$

From (43) we observe that the posterior variance of the random effects u_i is the sum of $d_{c_i}^{*-2}$ the inverse of the unique value of d_i^2 (the hyper-parameter of u_i), and all the elements in $\boldsymbol{\Sigma}_{\boldsymbol{\eta}_i}$. As we allow that $\boldsymbol{\Sigma}_{\boldsymbol{\eta}_i} \neq \boldsymbol{\Sigma}_{\boldsymbol{\eta}_j}$ ($\forall i \neq j$), s_i^2 is in turn allowed to be different for each individual effect u_i .

The covariance matrix of each composite error vector $\boldsymbol{\varepsilon}_i$ is also allowed to be different for every i . The covariance matrix of $\boldsymbol{\varepsilon}_i$ is given by

$$\text{Cov}(\boldsymbol{\varepsilon}_i) = \boldsymbol{\Sigma}_i = \boldsymbol{\Sigma}_{\boldsymbol{\eta}_i} + s_i^2 \iota_T \iota_T'. \quad (44)$$

4.2 MCMC Algorithm

For the choice of base distributions, we use the inverse gamma distribution, the conjugate prior for the variance of a normal distribution, i.e.

$$\begin{aligned} F_0 &\equiv \mathcal{IG}(a_u, b_u) \\ G_0 &\equiv \mathcal{IG}(a_\eta, b_\eta), \end{aligned} \quad (45)$$

where a_u and a_η are the shape hyper-parameters, and b_u and b_η denote, respectively, the rate hyper-parameters of F_0 and G_0 .

The likelihood of $\boldsymbol{\beta}$ marginalized over u_i is given by

$$p(\mathbf{y}_i | \boldsymbol{\beta}, \boldsymbol{\Sigma}_i) = \frac{1}{(2\pi)^{Q/2}} |\boldsymbol{\Sigma}_i|^{-\frac{1}{2}} \exp \left[-\frac{1}{2} (\mathbf{y}_i - \mathbf{X}_i \boldsymbol{\beta})' \boldsymbol{\Sigma}_i^{-1} (\mathbf{y}_i - \mathbf{X}_i \boldsymbol{\beta}) \right]. \quad (46)$$

²¹ c_{it} (c_{is}) denotes the group id's of η_{it} (η_{is}).

²²Here we adopt a prior whose mean is 0, and variance is 1000.

Compared with the marginal likelihood of the parametric Bayesian REM in (35), the covariance matrix of the composite error vector $\boldsymbol{\varepsilon}_i$ is allowed to be different for each individual i in the panel.

Given a conjugate normal prior for $\boldsymbol{\beta}$, i.e.

$$\boldsymbol{\beta} \sim \mathcal{N}(\mathbf{b}_0, \mathbf{V}_0),$$

where \mathbf{b}_0 and \mathbf{V}_0 denote respectively, prior mean and covariance matrix of $\boldsymbol{\beta}$, the posterior of $\boldsymbol{\beta}$ marginalized over u_i is

$$\boldsymbol{\beta} | \mathbf{y}, d_{c_i}^{*2}, \sigma_{c_{it}}^{*2} \sim \mathcal{N}(\mathbf{b}, \mathbf{V}). \quad (47)$$

$$\mathbf{V} = \left(\mathbf{V}_0^{-1} + \sum_{i=1}^N \mathbf{X}_i' \boldsymbol{\Sigma}_i^{-1} \mathbf{X}_i \right)^{-1}, \quad (48)$$

denotes the posterior covariance matrix, and \mathbf{b} is the posterior mean vector, which we write as

$$\mathbf{b} = \mathbf{V} \left(\mathbf{V}_0^{-1} \mathbf{b}_0 + \sum_{i=1}^N \mathbf{X}_i' \boldsymbol{\Sigma}_i^{-1} \mathbf{y}_i \right). \quad (49)$$

For $\Theta_\eta = \{\sigma_{it}^2\}$, $U = \{u_i\}$, and $\Theta_u = \{d_i^2\}$, a Gibbs sampler for this DP-REM can be written as:

$$\begin{aligned} & \Theta_\eta | \mathbf{y}, \mathbf{X}, U, \boldsymbol{\beta}, \Theta_u, \alpha_u, \alpha_\eta \\ & \Theta_u | \mathbf{y}, \mathbf{X}, U, \boldsymbol{\beta}, \Theta_\eta, \alpha_u, \alpha_\eta \\ & U | \mathbf{y}, \mathbf{X}, \boldsymbol{\beta}, \Theta_u, \Theta_\eta, \alpha_u, \alpha_\eta \\ & \boldsymbol{\beta} | \mathbf{y}, \mathbf{X}, U, \Theta_u, \Theta_\eta, \alpha_u, \alpha_\eta \\ & \alpha_u | \mathbf{y}, \mathbf{X}, U, \boldsymbol{\beta}, \Theta_u, \Theta_\eta, \alpha_\eta \\ & \alpha_\eta | \mathbf{y}, \mathbf{X}, U, \boldsymbol{\beta}, \Theta_u, \Theta_\eta, \alpha_u. \end{aligned} \quad (50)$$

The Gibbs sampler for the regression parameters, hyper-parameters and concentration parameters of the two \mathcal{DP} are similar to those for the DP-SUR in Section 3.2. In the DP-REM, the random effects have a mixture of normal distributions. The posterior mean and variance of each particular u_i are in (42) and (43), respectively. For each i a u_i is drawn from $\mathcal{N}(\mu_i, s_i^2)$ with the Gibbs sampler.

4.3 Correlated Random Effects Model

The Correlated Random Effects Model (CREM) represents a natural extension of the REM. Introduced by Mundlak (1978) and further discussed by Chamberlain (1980), the CREM offers a middle ground between the fixed and random effects.

Without loss of generality, we consider the following model for the panel data

$$y_{it} = \beta_1 x_{1it} + \beta_2 x_{2it} + v_i + \eta_{it}, \quad (51)$$

where v_i is the random effects. While maintaining the GLS structure of the REM, CREM allows the individual effects to be correlated with \mathbf{X}_i , representing the correlation using a linear function of the means of \mathbf{X}_i , i.e.

$$v_i = \beta_3 \bar{x}_{1i} + \beta_4 \bar{x}_{2i} + u_i, \quad (52)$$

The CREM model is then

$$y_{it} = \beta_1 x_{1it} + \beta_2 x_{2it} + \beta_3 \bar{x}_{1i} + \beta_4 \bar{x}_{2i} + u_i + \eta_{it} \quad (53)$$

The \mathcal{DP} prior can be introduced on the hyper-parameters of u_i and η_{it} as in the REM case.

4.4 DP-REM/CREM Simulation Results

We carry out a series of simulation experiments to demonstrate the performance of our DP-REM and DP-CREM methods relative to the standard Bayesian REM and CREM. The simulation experiments have been designed for the same purpose as those for the DP-SUR in Section 3.3.

For the REM model we assume

$$y_{it} = \beta_1 x_{1it} + \beta_2 x_{2it} + u_i + \eta_{it} = \beta_1 x_{1it} + \beta_2 x_{2it} + \varepsilon_{it} \quad (54)$$

where the explanatory variables are generated from the following normal distributions

$$x_{1,it} \stackrel{iid}{\sim} \mathcal{N}(1, 1), \quad x_{2,it} \stackrel{iid}{\sim} \mathcal{N}(3, 1).$$

We set the coefficients in (54) to

$$\beta_1 = 5, \quad \beta_2 = 10.$$

The coefficients in the CREM model (53) are set to

$$\beta_1 = 5, \quad \beta_2 = 10, \quad \beta_3 = -2, \quad \beta_4 = 2. \quad (55)$$

Below we present the simulation results. We first report the results where the errors, u_i and η_{it} , are assumed to follow t-distributions and then those with log-normal distributions.

t-Distributed Random Effects and Errors

Table 5 reports the averages of the posterior S.D.s of the REM coefficients estimated with both methods, and the average of the percentage differences between the DP-REM and REM posterior S.D.s. The largest differences between the two estimators with respect to the posterior S.D. are observed when $df = 2$, where the t-distributions of the random effects and the errors have the heaviest tails. As expected, these differences decrease as the df increase, where the tails of the t distributions become less 'heavy'. In the bottom panel where the errors have normal distributions (equivalent to df being infinity), the DP-REM and normal REM posterior S.D. are almost equivalent, as the t-distribution is the normal distribution in this case.

We also note that the percentage differences increase slightly when the sample size becomes larger for all three finite df . This is expected given that there are more extreme realizations in larger samples, and our DP-REM method detects such heterogeneity and assign them into the same group. In contrast, the Bayesian REM method assuming normality flattens the normal posterior distribution for the extreme values, leading to larger posterior S.D.

Table 6 reports the averages of posterior S.D. of the CREM coefficients, and the averages percentage differences between the two estimators. β_1 and β_2 denote the two original explanatory variables, whereas β_3 and β_4 capture the effect of the respective sample means for each individual in the panel. The findings are similar to the REM case in that the percentage differences between the posterior S.D. estimated with our DP-CREM and the parametric Bayesian CREM are the largest with df equal to 2, and decrease with the increase in the df . Differences between the two methods regarding the posterior S.D. are almost zero when the df is infinity, when the t-distribution becomes normal distribution. The percentage differences also increase slightly in the three finite df cases when the sample size becomes large due to more extreme values in the unobservables.

Table 5: Posterior s.d., REM with t distributed unobservables

df = 2									
Sample size	100			300			500		
Parameters	DP	REM	$\Delta\%$	DP	REM	$\Delta\%$	DP	REM	$\Delta\%$
β_1	0.0751	0.1419	43.18%	0.0484	0.0956	47.88%	0.0340	0.0671	47.92%
β_2	0.0499	0.0883	41.45%	0.0316	0.0575	47.04%	0.0213	0.0422	48.50%
df = 3									
Sample size	100			300			500		
Parameters	DP	REM	$\Delta\%$	DP	REM	$\Delta\%$	DP	REM	$\Delta\%$
β_1	0.0627	0.0803	20.61%	0.0366	0.0475	22.36%	0.0290	0.0379	22.94%
β_2	0.0419	0.0530	20.05%	0.0237	0.0309	22.65%	0.0183	0.0242	24.12%
df = 4									
Sample size	100			300			500		
Parameters	DP	REM	$\Delta\%$	DP	REM	$\Delta\%$	DP	REM	$\Delta\%$
β_1	0.0588	0.0667	11.51%	0.0346	0.0393	11.54%	0.0270	0.0310	12.59%
β_2	0.0394	0.0442	10.66%	0.0225	0.0257	12.15%	0.0171	0.0198	13.67%
df = ∞									
Sample size	100			300			500		
Parameters	DP	REM	$\Delta\%$	DP	REM	$\Delta\%$	DP	REM	$\Delta\%$
β_1	0.0484	0.0484	-0.02%	0.0274	0.0274	-0.20%	0.0211	0.0210	-0.19%
β_2	0.0301	0.0301	0.03%	0.0180	0.0179	-0.45%	0.0139	0.0139	0.01%

Log-normal Distributed Random Effects and Errors

Table 7 contains the average of posterior s.d. estimated with the DP-REM and normal REM, and the percentage differences between them. It can be seen that our DP-REM posterior s.d. are smaller than those estimated by Bayesian REM assuming normality in all cases. Due to the fact that the log-normal distribution is heavy tailed, the percentage differences are more than 70% in all cases, which increase slightly when the sample size gets larger.

The posterior s.d. of the DP-CREM and CREM averaged over the simulated samples are reported in Table 8, along with the average of the percentage difference between the two s.d. ASs before, the posterior s.d. estimated with our DP-CREM are more than 70% smaller than those estimated with the normal Bayesian CREM for all coefficients. The percentage differences also increase when the sample size increases

4.5 DP-REM/CREM Empirical Examples

In this section we present the results based upon two empirical examples. In the first we estimate the cost function of U.S. banks, and in the second we estimate the wages of U.S. workers.

Bank Cost Function

We first apply our DP-REM and DP-CREM methods to the dataset in Feng and Serletis (2009) on the costs of 218 U.S. banks whose assets are between 1 and 3 billion dollars (2000 value), covering a period of 8 years from 1998 to 2005. There are three inputs, labour, borrowed funds and physical capital; and three outputs, consumer loans, non-consumer loans and securities. The functional form is the simple translog cost function (Christensen and Greene, 1976). For

Table 6: Posterior s.d., CREM with t distributed unobservables

df = 2									
Sample size	100			300			500		
Parameters	DP	REM	$\Delta\%$	DP	REM	$\Delta\%$	DP	REM	$\Delta\%$
β_1	0.0860	0.1482	41.82%	0.0446	0.0935	48.33%	0.0394	0.0765	49.06%
β_2	0.0793	0.1409	41.60%	0.0471	0.0978	48.84%	0.0371	0.0790	50.10%
β_3	0.3665	0.6575	37.62%	0.2743	0.4854	42.63%	0.1595	0.3255	48.31%
β_4	0.1545	0.2768	40.26%	0.1197	0.2001	43.61%	0.0671	0.1412	49.43%
df = 3									
Sample size	100			300			500		
Parameters	DP	REM	$\Delta\%$	DP	REM	$\Delta\%$	DP	REM	$\Delta\%$
β_1	0.0689	0.0887	21.09%	0.0368	0.0472	21.51%	0.0294	0.0387	23.24%
β_2	0.0659	0.0845	20.99%	0.0383	0.0493	21.96%	0.0304	0.0398	22.92%
β_3	0.3077	0.3884	18.53%	0.1833	0.2434	23.10%	0.1373	0.1825	24.09%
β_4	0.1279	0.1619	19.54%	0.0749	0.0989	23.23%	0.0575	0.0763	24.06%
df = 4									
Sample size	100			300			500		
Parameters	DP	REM	$\Delta\%$	DP	REM	$\Delta\%$	DP	REM	$\Delta\%$
β_1	0.0639	0.0723	11.29%	0.0345	0.0399	12.97%	0.0256	0.0299	14.09%
β_2	0.0607	0.0689	11.60%	0.0358	0.0414	13.09%	0.0263	0.0306	14.04%
β_3	0.2857	0.3175	9.35%	0.1740	0.1982	11.91%	0.2547	0.2944	12.80%
β_4	0.1177	0.1316	10.09%	0.0707	0.0808	12.25%	0.0915	0.1056	12.78%
df = ∞									
Sample size	100			300			500		
Parameters	DP	REM	$\Delta\%$	DP	REM	$\Delta\%$	DP	REM	$\Delta\%$
β_1	0.0518	0.0517	-0.31%	0.0295	0.0295	-0.14%	0.0228	0.0228	-0.02%
β_2	0.0503	0.0502	-0.25%	0.0298	0.0296	-0.59%	0.0225	0.0224	-0.59%
β_3	0.2427	0.2416	-0.51%	0.1359	0.1365	0.34%	0.1076	0.1075	-0.10%
β_4	0.0952	0.0951	-0.13%	0.0586	0.0586	0.03%	0.0431	0.0431	-0.26%

Table 7: Posterior s.d., REM with log-normal distributed unobservables

Log-normal									
Sample size	100			300			500		
Parameters	DP	REM	$\Delta\%$	DP	REM	$\Delta\%$	DP	REM	$\Delta\%$
β_1	0.0733	0.3656	79.56%	0.0384	0.2228	82.43%	0.0305	0.1788	82.86%
β_2	0.0608	0.2350	73.37%	0.0311	0.1471	78.75%	0.0240	0.1163	79.02%

Table 8: Posterior s.d., CREM with log-normal distributed unobservables

Log-normal									
Sample size	100			300			500		
Parameters	DP	REM	$\Delta\%$	DP	REM	$\Delta\%$	DP	REM	$\Delta\%$
β_1	0.0754	0.3920	81.45%	0.0381	0.2479	85.24%	0.026	0.199	86.67%
β_2	0.0769	0.3785	81.37%	0.0390	0.2578	85.15%	0.024	0.194	87.17%
β_3	0.7722	1.7634	67.37%	0.3002	1.1358	74.96%	0.157	0.898	80.77%
β_4	0.3342	0.7583	70.91%	0.1080	0.4828	78.79%	0.063	0.378	82.24%

industry i with n inputs and m outputs we write

$$\begin{aligned} \ln C_{it} = & \sum_{j=1}^m \alpha_j \ln q_{j,it} + \frac{1}{2} \sum_{j=1}^m \sum_{k=1}^m \delta_{jk} \ln q_{j,it} \cdot \ln q_{k,it} + \sum_{r=1}^n \beta_r \ln p_{r,it} \\ & + \frac{1}{2} \sum_{r=1}^n \sum_{s=1}^n \phi_{rs} \ln p_{r,it} \cdot \ln p_{s,it} + \sum_{r=1}^n \sum_{j=1}^m \gamma_{rj} \ln p_{r,it} \cdot \ln q_{j,it} + u_i + \eta_{it}. \end{aligned} \quad (56)$$

where C is cost, q_j is output quantity j , and p_r is input price r . We impose the linear homogeneity in input prices on the cost function, which in the translog case can be expressed as

$$\begin{aligned} \sum_{r=1}^n \beta_r &= 1, \\ \sum_{s=1}^n \phi_{sr} &= 0, \quad r = 1, 2, \dots, n, \\ \sum_{r=1}^n \gamma_{rj} &= 0, \quad j = 1, 2, \dots, m. \end{aligned} \quad (57)$$

Table 9 contains the posterior means and s.d. of the free coefficients in the REM. To differentiate the inputs from outputs, we index the three outputs with numbers, and index the inputs by letters, with l and f denoting, respectively, labour and borrowed funds. The posterior mode of the number of groups²³ in the random effects is 2, and that in the errors is 3, which shows the existence of heterogeneity, although it is not strong.

The posterior means of all the coefficients for first order terms (the β 's and α 's) are of the same magnitude with the two methods with the exception of α_1 for customer loans, but it is insignificant with the REM. Although a number of the coefficients for the crossproduct terms have different signs across the two methods, the coefficients are not significant. Consistent with the detection of heterogeneity in the random effects and the errors, the posterior s.d. of our DP-SUR method are smaller than those estimated by parametric Bayesian REM for all coefficients. Most of the percentage differences presented in the last column are more than 10%, with the largest being 24.51% for δ_{33} .

We also estimate the model with the DP-CREM. The posterior mode of the number of groups in the random effects and the errors are 2 and 3, respectively, indicating the existence of heterogeneity. The coefficients of the explanatory variables have smaller DP-CREM posterior s.d. than their CREM counterparts, with similar magnitudes as in Table 9. However, the regression parameters are all highly insignificant for the sample means of the explanatory variables, as their posterior s.d. are very large compared with their posterior means. This indicates that the data do not support the CREM specification²⁴.

U.S. Individual Wage

In this section we present the results of a wage model for U.S. workers using the data in Cornwell and Rupert (1988). The data covers 595 individuals over a period of 7 years, from 1976 to 1982. This sample size allows us to demonstrate our method with sub-samples of 100 and 250 individuals.

²³The number of groups, i.e. the number of normal distributions that are being mixed, is also a random variable in Bayesian methods. Its posterior mode thus provides an indication of the strength of heterogeneity in the error distributions.

²⁴The results are not presented here, but are available on request.

Table 9: U.S. Bank Cost Function REM

Coefficients	Mean		S.D.		
	DP-REM	REM	DP-REM	REM	$\Delta\%$
β_l	-0.6583	-0.6099	0.1480	0.1795	17.56%
β_f	1.4475	1.0758	0.1068	0.1091	2.12%
α_1	-0.1908	-0.0240	0.0799	0.0944	15.35%
α_2	0.6870	0.5777	0.1405	0.1562	10.05%
α_3	0.9486	0.8439	0.1284	0.1558	17.55%
ϕ_{ll}	0.1087	0.0405	0.0170	0.0214	20.44%
ϕ_{ff}	0.1670	0.1266	0.0055	0.0071	22.37%
ϕ_{lf}	-0.1664	-0.0891	0.0079	0.0097	18.57%
δ_{11}	0.0141	0.0145	0.0033	0.0041	20.02%
δ_{22}	0.1267	0.1098	0.0192	0.0218	12.13%
δ_{33}	0.0934	0.1030	0.0134	0.0178	24.51%
δ_{12}	0.0033	-0.0113	0.0067	0.0082	17.41%
δ_{13}	0.0011	-0.0032	0.0048	0.0063	23.39%
δ_{23}	-0.1441	-0.1318	0.0117	0.0152	22.98%
γ_{l1}	0.0049	0.0181	0.0057	0.0063	9.82%
γ_{l2}	-0.0030	0.0370	0.0131	0.0151	13.60%
γ_{l3}	0.0086	-0.0030	0.0112	0.0135	17.36%
γ_{f1}	-0.0157	-0.0220	0.0035	0.0043	19.51%
γ_{f2}	0.0172	-0.0241	0.0088	0.0100	12.21%
γ_{f3}	0.0217	0.0570	0.0064	0.0080	19.82%

The model is given by

$$\ln Wage_{it} = \beta_1 E_{it} + \beta_2 M_{it} + \beta_3 F_i + \beta_4 Ed_{it} + v_i + \varepsilon_{it},$$

where the dependent variable is the logged wage, and the explanatory variables are experience in years (E), dummies for marriage status (M) and the individual being female (F), as well as the years of education (Ed). As there is a strong reason to suspect that the unobserved individual effect v_i to be endogenous due to omitted variables such as personal capability and motivation, we apply our DP-CREM model, and write v_i as

$$v_i = \tilde{\beta}_1 \bar{E}_i + \tilde{\beta}_2 \bar{M}_i + u_i. \quad (58)$$

The means of experience and marriage status of individual i are included as they are the two time variant variables in the original model.

Table 10: U.S. Individual Wage CREM

Sample size	Mean				S.D.					
	100		250		100			250		
Parameters	DP	REM	DP	REM	DP	REM	$\Delta\%$	DP	REM	$\Delta\%$
β_1	0.100	0.102	0.094	0.099	0.0025	0.0032	24.18%	0.0015	0.0020	25.49%
β_2	0.027	0.038	-0.043	-0.070	0.0375	0.0524	28.43%	0.0231	0.0308	25.07%
β_3	5.184	1.945	5.077	1.825	0.2646	0.4256	37.83%	0.1576	0.2602	39.42%
β_4	0.085	0.334	0.075	0.288	0.0189	0.0213	11.22%	0.0091	0.0131	30.74%
$\tilde{\beta}_1$	-0.091	-0.058	-0.085	-0.058	0.0046	0.0077	40.21%	0.0028	0.0056	50.17%
$\tilde{\beta}_2$	5.450	1.603	5.668	2.261	0.2968	0.3381	12.23%	0.1493	0.2091	28.60%

Table 10 contains the posterior means, S.D. and the percentage difference between the S.D. estimated with our DP-CREM and Bayesian CREM assuming normality for both the 100 and 250 individual sub-sample. The two coefficients for the means of time variant explanatory variables, $\tilde{\beta}_1$ and $\tilde{\beta}_2$ are both significant with both sub-samples, indicating that v_i actually is correlated to

the explanatory variables, confirming our suspicion. As for the coefficients for the explanatory variables themselves (β_1 to β_4), the posterior means are all of the same signs with the DP-CREM and DP-REM, and the differences between them become smaller with the larger sub-sample of 250 observations. As expected, the experience and education are both positively correlated with the wage of the workers. The coefficient for the gender dummy (β_3) is also positive. This may seem as an indication of gender discrimination in wages against male workers within the two sub-samples²⁵. Heterogeneity is detected in both sub-samples, as the posterior modes of the numbers of groups in the random effects and errors are 2 and 3 with the 100 individual sub-sample, and 3 and 4 with the 250 individual one. Our semi-parametric DP-CREM provides smaller posterior S.D. for all coefficients with both sub-samples.

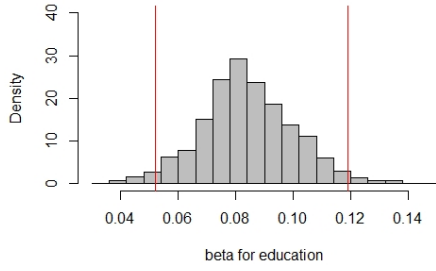
In Figure 2 we present histograms of the posterior distributions of the education (β_4) parameter. As before, the two red vertical lines in each panel mark the 95% credible interval. Comparing panel (a) and (b) we observe that for the 100-observation sub-sample both the DP-CREM and the parametric Bayesian CREM give 95% credible that exclude 0, with the 95% credible interval for DP-CREM (from 0.052 to 0.119) shorter than the parametric CREM (from 0.292 to 0.377). This is not surprising given that the posterior distribution is right skewed as shown in panel (a). Similar conclusions follow from comparing panel (c) and (d) based upon the 250-observation sub-sample. The 95% credible interval with the DP-CREM is from 0.058 to 0.092, while that with the parametric CREM is from 0.263 to 0.313. Note also that the 95% credible intervals based upon the 250-observation sub-sample are shorter than their counterparts using the 100-observation sub-sample, due to the increase in the sample size.

5 Conclusion

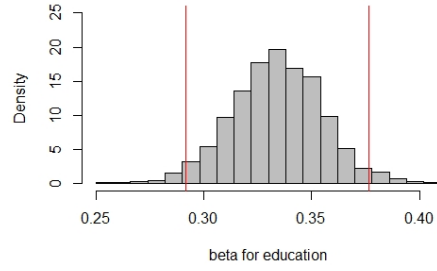
In this paper we address the potential violation of the assumptions made by parametric Bayesian GLS estimators that the unobservables are homogeneous regarding their distributions. Such assumptions are likely to be problematic in reality particularly when micro data are used, as the features of individuals or households are likely to lead to the distributions of observations being different. We present a semi-parametric Bayesian GLS where the error distribution is a non-parametric mixture of normal distributions by introducing a Dirichlet process prior on the hyper-parameters of the errors. The number of normal components is decided jointly by the data and the prior in such a mixture of normals, which is able to cover a large variety of distributions. The errors are grouped by the \mathcal{DP} prior, with those in the same group having the same hyper-parameters and thus the same distribution. Two specific cases of the semi-parametric Bayesian GLS are then introduced, which are the SUR for equation systems and the REM/CREM for panel data.

Our DP-SUR and DP-REM/DP-CREM methods are demonstrated with a series of simulation experiments consisting of three scenarios, where the unobservables follow normal distributions, t-distributions which are one type of scale mixtures of normals, and log-normal distributions, respectively. The results show that in the homogeneous normal case, our DP-SUR and DP-REM/DP-CREM methods give posterior means and S.D. similar to their parametric counterparts assuming normality. When the errors follow t-distributions, the degrees of freedom of the t-distribution control how heavy the tails are, which reflects the strength of heterogeneity in the unobservables. Our simulation results show that the posterior S.D. of our DP-SUR and

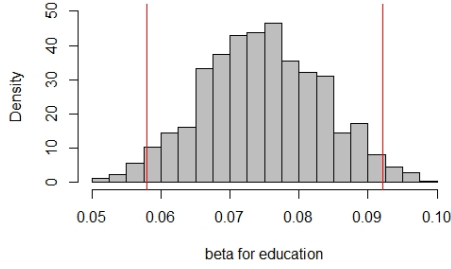
²⁵However, before jumping to this conclusion, one should keep in mind that in the time period of the data (1976 to 1982), fewer women were working than at present. In the 100-observation sub-sample, 16% of the workers are women, and in the 250-observation sub-sample, the percentage is 8.8%. Given such low percentage of female workers, those women who did decide to enter the labour market may themselves be relatively skilled workers who could reasonably expect a high wage. The sample selection bias (Heckman, 1976) may still be present here as a result, though it is not within the scope of this paper.



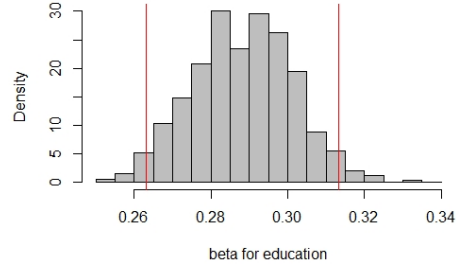
(a) DP-CREM, 100-observations



(b) Parametric CREM, 100-observations



(c) DP-CREM, 250-observations



(d) Parametric CREM, 250-observations

Figure 2: Histograms of the parameter for education (β_4)

DP-REM/DP-CREM are smaller than those of the parametric Bayesian methods. Such efficiency gains are the largest when df is 2 that represents the strongest heterogeneity. The efficiency gains become smaller when the df increases, for the tails are less heavy, i.e. the heterogeneity is less strong. The simulations with log-normal unobservables are used to demonstrate the robustness of our method with asymmetric, fat tailed distributions. The results demonstrated that the posterior S.D. of our DP-SUR and DP-REM/DP-CREM method are more than 50% smaller than those of the parametric Bayesian estimators assuming normality. Moreover, the efficiency gains increase slightly with larger sample sizes when the distribution of the unobservables are non-normal, which is the result of more extreme realizations in large samples.

We apply our DP-SUR method to the demands for production factors with the generalized Leontief cost function using a dataset of the U.S. banking industry. We estimate the model with an 800-observation sub-sample and the full sample. Heterogeneity is detected in both the sub-sample and the full sample. The DP-SUR posterior S.D. are smaller than the normal Bayesian SUR ones for all the demand elasticities, which shows that it is more preferable to use a semi-parametric method such as our DP-SUR.

Our DP-REM/DP-CREM are applied to two datasets as well. The first is a U.S. bank cost functions data. The REM seems to fit the datasets better. Heterogeneity is detected in the U.S. bank data, and our DP-REM achieved smaller posterior S.D. than the parametric Bayesian REM. The second application is a U.S. individual wage model, where there is a strong reason to suspect that the unobserved individual effects are correlated to the explanatory variables like education due to unobserved individual features such as abilities. The CREM model is then estimated. Our DP-CREM detects heterogeneity in this dataset as well, and obtains smaller posterior S.D. than the parametric Bayesian CREM.

A Posterior Means for DP-SUR Simulations

A.1 Multivariate t-distributed Errors

Table 11 gives the posterior means averaged over the samples that are estimated with the DP-SUR and SUR assuming normality. One can see that the posterior means estimated with both our semi-parametric DP-SUR and the Bayesian SUR assuming normality are similar to each other. In addition, they are both close to the true values of the coefficients in all cases.

Table 11: Posterior means, multivariate t errors

df = 2							
Sample size	Truth	100		250		500	
Parameters		DP	SUR	DP	SUR	DP	SUR
β_{10}	1	1.035	0.942	1.017	1.044	0.986	1.051
β_{11}	-0.5	-0.500	-0.484	-0.491	-0.488	-0.500	-0.519
β_{12}	1.6	1.590	1.606	1.591	1.585	1.606	1.590
β_{20}	1.5	1.524	1.568	1.482	1.540	1.484	1.387
β_{21}	-1.2	-1.195	-1.194	-1.201	-1.194	-1.206	-1.230
β_{22}	-0.7	-0.702	-0.719	-0.697	-0.700	-0.696	-0.684
β_{23}	2	2.003	2.006	1.998	1.980	2.010	2.004
df = 3							
Sample size	Truth	100		250		500	
Parameters		DP	SUR	DP	SUR	DP	SUR
β_{10}	1	1.008	1.002	1.015	1.030	0.999	0.989
β_{11}	-0.5	-0.497	-0.493	-0.502	-0.499	-0.498	-0.499
β_{12}	1.6	1.600	1.600	1.596	1.591	1.598	1.601
β_{20}	1.5	1.507	1.507	1.493	1.502	1.503	1.504
β_{21}	-1.2	-1.194	-1.183	-1.198	-1.199	-1.202	-1.199
β_{22}	-0.7	-0.699	-0.696	-0.698	-0.699	-0.703	-0.702
β_{23}	2	1.995	1.989	1.999	2.001	1.995	1.996
df = 4							
Sample size	Truth	100		250		500	
Parameters		DP	SUR	DP	SUR	DP	SUR
β_{10}	1	0.973	0.979	1.000	1.000	0.989	0.992
β_{11}	-0.5	-0.500	-0.499	-0.516	-0.519	-0.495	-0.495
β_{12}	1.6	1.610	1.609	1.606	1.608	1.603	1.600
β_{20}	1.5	1.516	1.536	1.519	1.522	1.509	1.495
β_{21}	-1.2	-1.205	-1.208	-1.192	-1.192	-1.202	-1.202
β_{22}	-0.7	-0.708	-0.713	-0.700	-0.701	-0.704	-0.701
β_{23}	2	1.996	1.997	2.000	1.995	1.998	1.999
df = ∞							
Sample size	Truth	100		250		500	
Parameters		DP	SUR	DP	SUR	DP	SUR
β_{10}	1	0.9895	0.9864	0.9853	0.9846	0.9813	0.9835
β_{11}	-0.5	-0.4969	-0.4974	-0.4959	-0.4957	-0.4904	-0.4905
β_{12}	1.6	1.6030	1.6040	1.6040	1.6044	1.6005	1.5999
β_{20}	1.5	1.4706	1.4706	1.5122	1.5122	1.4675	1.4677
β_{21}	-1.2	-1.2115	-1.2118	-1.1959	-1.1965	-1.1976	-1.1977
β_{22}	-0.7	-0.6984	-0.6985	-0.7004	-0.7003	-0.6911	-0.6912
β_{23}	2	1.9950	1.9953	2.0000	2.0004	1.9973	1.9979

A.2 Multivariate Log-normal Errors

Table 12 contains the posterior means estimated with the two methods with multivariate log-normal errors. In all three samples the two posterior means of all the slope parameters are similar, and are close to the truth. The intercepts β_{10} and β_{20} estimated with our DP-SUR, however, are farther away from the true values. The fact that the log-normal distribution is skewed influences the posterior means of intercepts, when the \mathcal{DP} mixture model mixes normal distributions to model the log-normal distribution.

Table 12: Posterior means, multivariate log-normal errors

		Log-normal					
Sample size	Truth	100		250		500	
Parameters		DP	SUR	DP	SUR	DP	SUR
β_{10}	1	0.199	0.996	0.130	0.973	0.116	0.987
β_{11}	-0.5	-0.501	-0.492	-0.499	-0.488	-0.496	-0.504
β_{12}	1.6	1.600	1.598	1.606	1.605	1.602	1.603
β_{20}	1.5	0.697	1.553	0.665	1.557	0.654	1.568
β_{21}	-1.2	-1.201	-1.180	-1.202	-1.205	-1.199	-1.198
β_{22}	-0.7	-0.701	-0.696	-0.703	-0.719	-0.705	-0.720
β_{23}	2	2.004	2.012	2.002	1.993	1.995	1.988

B Posterior Standard Deviations for DP-SUR Simulations

B.1 Multivariate t-distributed Errors

Table 13 presents the full results regarding the posterior standard deviations of the DP-SUR simulations with multivariate t-distributed errors.

B.2 Multivariate Log-normal Errors

Table 14 shows the full results regarding the posterior standard deviations of the DP-SUR simulations with multivariate log-normal errors.

C Posterior Means for DP-REM/CREM

C.1 t-distributed Errors

Table 15 contains the average of the posterior means over the samples of the coefficients in REM with t-distributed random effects and errors. It can be seen that the averaged posterior means estimated with our DP-REM and parametric Bayesian REM are all almost identical, and they are all close to the true value of the coefficients for all four cases, i.e. df being 2, 3, 4, and infinity, where the t-distribution becomes normal distribution.

The average of posterior means of the CREM coefficients are presented in Table 16. Similar to the REM case, the average of the posterior means estimated with both our DP-CREM and parametric Bayesian CREM are similar to each other, and close to the pre-set true values of the coefficients in all cases.

C.2 Log-normal Errors

Table 17 gives the average of the posterior means of the REM with log-normal distributed random effects and errors. One could see that the the DP-REM and Bayesian REM assuming normality

Table 13: Posterior s.d., multivariate t errors, full results

df = 2									
Sample size	100			250			500		
Parameters	DP	SUR	$\Delta\%$	DP	SUR	$\Delta\%$	DP	SUR	$\Delta\%$
β_{10}	0.4184	0.8103	42.59%	0.2328	0.5262	50.09%	0.1644	0.4175	55.49%
β_{11}	0.1149	0.2155	41.57%	0.0732	0.1596	49.56%	0.0458	0.1126	54.22%
β_{12}	0.1254	0.2358	41.74%	0.0696	0.1515	49.30%	0.0498	0.1223	54.12%
β_{20}	0.5085	0.9791	42.67%	0.3261	0.7255	48.60%	0.2366	0.5649	53.30%
β_{21}	0.1107	0.2080	41.49%	0.0649	0.1359	48.48%	0.0508	0.1198	52.78%
β_{22}	0.1136	0.2141	41.63%	0.0653	0.1366	48.28%	0.0487	0.1139	52.34%
β_{23}	0.1271	0.2391	41.63%	0.0631	0.1321	48.38%	0.0467	0.1098	52.55%
df = 3									
Sample size	100			250			500		
Parameters	DP	SUR	$\Delta\%$	DP	SUR	$\Delta\%$	DP	SUR	$\Delta\%$
β_{10}	0.4121	0.5327	20.63%	0.2250	0.3098	26.25%	0.1619	0.2289	28.09%
β_{11}	0.1114	0.1452	21.20%	0.0708	0.0973	26.05%	0.0446	0.0634	28.53%
β_{12}	0.1231	0.1586	20.40%	0.0670	0.0923	26.20%	0.0486	0.0687	28.01%
β_{20}	0.4980	0.6497	21.15%	0.3174	0.4344	25.73%	0.2260	0.3220	28.82%
β_{21}	0.1079	0.1397	20.63%	0.0630	0.0866	25.91%	0.0488	0.0694	28.70%
β_{22}	0.1114	0.1442	20.51%	0.0636	0.0869	25.63%	0.0462	0.0661	29.11%
β_{23}	0.1230	0.1609	21.35%	0.0617	0.0844	25.61%	0.0449	0.0638	28.70%
df = 4									
Sample size	100			250			500		
Parameters	DP	SUR	$\Delta\%$	DP	SUR	$\Delta\%$	DP	SUR	$\Delta\%$
β_{10}	0.3931	0.4570	13.98%	0.2211	0.2625	15.19%	0.1578	0.1913	17.09%
β_{11}	0.1079	0.1248	13.52%	0.0696	0.0826	15.22%	0.0440	0.0529	16.58%
β_{12}	0.1179	0.1363	13.49%	0.0659	0.0783	15.32%	0.0476	0.0574	16.81%
β_{20}	0.4856	0.5568	12.78%	0.3091	0.3686	15.49%	0.2210	0.2702	17.82%
β_{21}	0.1054	0.1200	12.15%	0.0618	0.0733	15.08%	0.0473	0.0582	18.41%
β_{22}	0.1086	0.1236	12.17%	0.0617	0.0739	15.88%	0.0456	0.0554	17.30%
β_{23}	0.1204	0.1376	12.56%	0.0603	0.0717	15.19%	0.0438	0.0534	17.77%
df = ∞									
Sample size	100			250			500		
Parameters	DP	SUR	$\Delta\%$	DP	SUR	$\Delta\%$	DP	SUR	$\Delta\%$
β_{10}	0.3076	0.3016	-2.04%	0.2055	0.2031	-1.20%	0.1275	0.1270	-0.39%
β_{11}	0.0854	0.0837	-2.09%	0.0569	0.0565	-0.85%	0.0372	0.0371	-0.56%
β_{12}	0.0909	0.0891	-1.97%	0.0611	0.0605	-1.13%	0.0366	0.0365	-0.42%
β_{20}	0.4796	0.4696	-2.19%	0.2754	0.2730	-0.93%	0.1826	0.1810	-0.96%
β_{21}	0.0879	0.0860	-2.25%	0.0575	0.0572	-0.58%	0.0388	0.0384	-1.14%
β_{22}	0.0974	0.0955	-2.03%	0.0568	0.0564	-0.83%	0.0395	0.0395	-0.25%
β_{23}	0.0888	0.0868	-2.31%	0.0557	0.0552	-0.93%	0.0355	0.0355	-0.23%

Table 14: Posterior s.d., multivariate log-normal errors, full results

Log-normal									
Sample size	100			250			500		
Parameters	DP	SUR	$\Delta\%$	DP	SUR	$\Delta\%$	DP	SUR	$\Delta\%$
β_{10}	0.2221	0.5537	56.72%	0.1536	0.4525	63.95%	0.1001	0.2927	64.94%
β_{11}	0.0720	0.1831	57.64%	0.0415	0.1266	65.06%	0.0272	0.0834	66.61%
β_{12}	0.0672	0.1741	58.37%	0.0438	0.1338	65.20%	0.0280	0.0867	66.95%
β_{20}	0.3363	0.8722	57.82%	0.2150	0.6088	63.42%	0.1382	0.4064	65.06%
β_{21}	0.0800	0.2119	58.78%	0.0439	0.1277	64.38%	0.0270	0.0826	66.39%
β_{22}	0.0716	0.1890	58.64%	0.0438	0.1277	64.44%	0.0297	0.0896	65.87%
β_{23}	0.0662	0.1749	58.58%	0.0398	0.1147	64.05%	0.0270	0.0806	65.62%

Table 15: Posterior means, REM with t distributed unobservables

df = 2							
Sample size	Truth	100		300		500	
Parameters		DP	REM	DP	REM	DP	REM
β_1	5	5.003	5.005	5.002	5.000	5.008	5.013
β_2	10	9.999	10.002	10.000	9.997	9.996	9.990
df = 3							
Sample size	Truth	100		300		500	
Parameters		DP	REM	DP	REM	DP	REM
β_1	5	4.997	4.996	5.004	5.005	4.999	4.999
β_2	10	10.001	10.000	9.998	9.997	9.999	9.999
df = 4							
Sample size	Truth	100		300		500	
Parameters		DP	REM	DP	REM	DP	REM
β_1	5	4.999	4.999	4.998	4.998	4.997	4.998
β_2	10	10.001	10.001	10.001	10.003	9.998	9.998
df = ∞							
Sample size	Truth	100		300		500	
Parameters		DP	REM	DP	REM	DP	REM
β_1	5	5.003	5.003	5.000	5.000	4.998	4.998
β_2	10	9.998	9.999	9.999	9.999	10.000	10.000

Table 16: Posterior means, CREM with t distributed unobservables

df = 2							
Sample size	Truth	100		300		500	
Parameters		DP	REM	DP	REM	DP	REM
β_1	5	4.999	5.003	4.996	4.988	5.001	4.997
β_2	10	10.002	9.996	10.001	10.005	10.000	9.999
β_3	-2	-1.973	-2.020	-1.988	-1.961	-2.005	-2.003
β_4	2	1.991	2.023	1.993	1.990	2.006	2.011
df = 3							
Sample size	Truth	100		300		500	
Parameters		DP	REM	DP	REM	DP	REM
β_1	5	4.996	4.993	5.002	5.005	4.994	4.995
β_2	10	10.000	9.997	10.000	10.003	9.997	9.995
β_3	-2	-1.987	-1.970	-1.983	-1.968	-1.986	-1.996
β_4	2	1.995	1.991	1.997	1.987	1.999	2.004
df = 4							
Sample size	Truth	100		300		500	
Parameters		DP	REM	DP	REM	DP	REM
β_1	5	4.998	4.999	5.002	5.001	4.999	4.997
β_2	10	10.003	10.003	10.002	10.001	9.998	9.993
β_3	-2	-1.999	-1.996	-1.999	-1.996	-2.005	-2.004
β_4	2	2.000	1.999	1.996	1.996	2.008	2.014
df = ∞							
Sample size	Truth	100		300		500	
Parameters		DP	REM	DP	REM	DP	REM
β_1	5	5.000	5.000	4.998	4.998	5.003	5.003
β_2	10	10.001	10.001	9.998	9.998	10.001	10.001
β_3	-2	-1.999	-2.000	-2.004	-2.005	-1.998	-1.998
β_4	2	1.998	1.998	2.004	2.005	1.999	1.999

obtain relatively close means to the true values of the coefficients, with the DP-REM posterior means being slightly closer to the truths.

Table 17: Posterior means, REM with log-normal distributed unobservables

Log-normal							
Sample size	Truth	100		300		500	
Parameters		DP	REM	DP	REM	DP	REM
β_1	5	5.147	5.491	5.124	5.414	5.113	5.371
β_2	10	10.357	11.163	10.369	11.232	10.359	11.222

Table 18 presents the posterior means averaged over the simulation samples of the DP-CREM and normal Bayesian CREM with log-normal distributed random effects and errors. The posterior means of the coefficients for the explanatory variables are very close to the truth with both our DP-CREM and Bayesian CREM assuming normality. The posterior means of the coefficients for the means of the explanatory variables are slightly farther away from the truth. The skewness of the log-normal distribution influences the posterior means of the intercepts, which are time invariant for each individual i in the random effects. In the CREM case, the sample means of each individual's explanatory variables, \bar{x}_{1i} and \bar{x}_{2i} , are also time invariant like the intercept. As a result, the posterior means of their coefficients, β_3 and β_4 , are more different from the truth compared with β_1 and β_2 , as the log-normal distribution is skewed.

Table 18: Posterior means, CREM with log-normal distributed unobservables

Log-normal							
Sample size	Truth	100		300		500	
Parameters		DP	REM	DP	REM	DP	REM
β_1	5	5.002	5.023	4.999	4.974	5.000	5.040
β_2	10	9.999	10.020	10.007	10.015	10.000	9.995
β_3	-2	-1.908	-1.269	-1.851	-1.657	-1.832	-1.432
β_4	2	2.669	3.734	2.590	3.896	2.577	3.818

References

- [1] Aldous, D. J. (1985). Exchangeability and related topics. In *École d'Été de Probabilités de Saint-Flour XIII—1983*. Springer, Berlin, 1-198.
- [2] Allenby, G. M., Arora, N., & Ginter, J. L. (1998). On the heterogeneity of demand. *Journal of Marketing Research*, 384-389.
- [3] Andrews, D. F., & Mallows, C. L. (1974). Scale mixtures of normal distributions. *Journal of the Royal Statistical Society. Series B (Methodological)*, 99-102.
- [4] Antoniak, C. E. (1974). Mixtures of Dirichlet processes with applications to Bayesian non-parametric problems. *The Annals of Statistics*, 1152-1174.
- [5] de Carvalho, V. I., Jara, A., Hanson, T. E., & de Carvalho, M. (2013). Bayesian nonparametric ROC regression modeling. *Bayesian Analysis*, 8(3), 623-646.
- [6] Chao, J. C., & Phillips, P. C. (1998). Posterior distributions in limited information analysis of the simultaneous equations model using the Jeffreys prior. *Journal of Econometrics*, 87(1), 49-86.
- [7] Chigira, H., & Shiba, T. (2015). Dirichlet Prior for Estimating Unknown Regression Error Heteroskedasticity. *TERG Discussion Papers*, 341, 1-17.
- [8] Clementi, F., & Gallegati, M. (2005). Pareto's law of income distribution: Evidence for Germany, the United Kingdom, and the United States. In *Econophysics of Wealth Distributions* (pp. 3-14). Springer, Milano.
- [9] Conley, T. G., Hansen, C. B., McCulloch, R. E., & Rossi, P. E. (2008). A semi-parametric Bayesian approach to the instrumental variable problem. *Journal of Econometrics*, 144(1), 276-305.
- [10] Cornwell, C., & Rupert, P. (1988). Efficient estimation with panel data: An empirical comparison of instrumental variables estimators. *Journal of Applied Econometrics*, 3(2), 149-155.
- [11] Christensen, L. R., & Greene, W. H. (1976). Economies of scale in US electric power generation. *Journal of political Economy*, 84(4, Part 1), 655-676.
- [12] Diewert, W. E. (1971). An application of the Shephard duality theorem: a generalized Leontief production function. *Journal of Political Economy*, 79(3), 481-507.
- [13] Escobar, M. D., & West, M. (1995). Bayesian density estimation and inference using mixtures. *Journal of the american statistical association*, 90(430), 577-588.
- [14] Escobar, M. D., & West, M. (1998). Computing nonparametric hierarchical models. In: Dey, D., Müller, P., & Sinha, D. *Practical nonparametric and semiparametric Bayesian statistics*. Springer, New York, 1-22.
- [15] Feng, G., & Serletis, A. (2009). Efficiency and productivity of the US banking industry, 1998 – 2005: Evidence from the Fourier cost function satisfying global regularity conditions. *Journal of Applied Econometrics*, 24(1), 105-138.
- [16] Ferguson, T. S. (1973). A Bayesian analysis of some nonparametric problems. *The Annals of Statistics*, 1(2), 209-230.

- [17] Gershman, S. J., & Blei, D. M. (2012). A tutorial on Bayesian nonparametric models. *Journal of Mathematical Psychology*, 56(1), 1-12.
- [18] Geweke, J. (1993). Bayesian treatment of the independent student-t linear model. *Journal of Applied Econometrics*, 8(S1).
- [19] Geweke, J. (1996). Bayesian reduced rank regression in econometrics. *Journal of Econometrics*, 75(1), 121-146.
- [20] Greene, William H. (2012). *Econometric Analysis*, Seventh Edition. Upper Saddle River, New Jersey: Prentice Hall.
- [21] Heckman, J. J. (1976). The common structure of statistical models of truncation, sample selection and limited dependent variables and a simple estimator for such models. *Annals of Economic and Social Measurement*, 5(4), 475-492. NBER.
- [22] Hejblum, B. P., Alkhassim, C., Gottardo, R., Caron, F., & Thiébaud, R. (2019). Sequential Dirichlet process mixtures of multivariate skew t -distributions for model-based clustering of flow cytometry data. *The Annals of Applied Statistics*, 13(1), 638-660.
- [23] Kleinman, K. P., & Ibrahim, J. G. (1998). A semiparametric Bayesian approach to the random effects model. *Biometrics*, 921-938.
- [24] Kleibergen, F., & van Dijk, H. K. (1998). Bayesian simultaneous equations analysis using reduced rank structures. *Econometric Theory*, 14(6), 701-743.
- [25] Kloek, T., & Van Dijk, H. K. (1978). Bayesian estimates of equation system parameters: an application of integration by Monte Carlo. *Econometrica*, 46, 1-19.
- [26] Koop, G. (2003). *Bayesian Econometrics*. John Wiley & Sons.
- [27] Kumbhakar, S. C., & Tsionas, E. G. (2011). Stochastic error specification in primal and dual production systems. *Journal of Applied Econometrics*, 26(2), 270-297.
- [28] Kyung, M., Gill, J., & Casella, G. (2010). Estimation in Dirichlet random effects models. *The Annals of Statistics*, 38(2), 979-1009.
- [29] Li, M., & Tobias, J. L. (2011). Bayesian inference in a correlated random coefficients model: Modeling causal effect heterogeneity with an application to heterogeneous returns to schooling. *Journal of econometrics*, 162(2), 345-361.
- [30] Li, C., Casella, G., & Ghosh, M. (2018). Estimation of regression vectors in linear mixed models with Dirichlet process random effects. *Communications in Statistics-Theory and Methods*, 47(16), 3935-3954.
- [31] MacEachern, S. N. (1998). Computational methods for mixture of Dirichlet process models. In: Dey, D., Müller, P., & Sinha, D. *Practical nonparametric and semiparametric Bayesian statistics*. Springer, New York, 23-43.
- [32] Malikov, E., Kumbhakar, S. C., & Tsionas, M. G. (2016). A cost system approach to the stochastic directional technology distance function with undesirable outputs: the case of US banks in 2001-2010. *Journal of Applied Econometrics*, 31(7), 1407-1429.
- [33] Murtazashvili, I., & Wooldridge, J. M. (2008). Fixed effects instrumental variables estimation in correlated random coefficient panel data models. *Journal of Econometrics*, 142(1), 539-552.

- [34] Rossi, P. E., Allenby, G. M., & McCulloch, R. (2012). *Bayesian statistics and marketing*. John Wiley & Sons.
- [35] Strutz, T. (2010). *Data fitting and uncertainty: A practical introduction to weighted least squares and beyond*. Vieweg and Teubner.
- [36] Teh, Y. W., Jordan, M. I., Beal, M. J., & Blei, D. M. (2005). Sharing clusters among related groups: Hierarchical Dirichlet processes. In *Advances in neural information processing systems* (pp. 1385-1392).
- [37] Teh, Y. W. (2011). Dirichlet Process. In *Encyclopedia of machine learning*, pp. 280-287. Springer US.
- [38] Wiesenfarth, M., Hisgen, C. M., Kneib, T., & Cadarso-Suarez, C. (2014). Bayesian non-parametric instrumental variables regression based on penalized splines and dirichlet process mixtures. *Journal of Business & Economic Statistics*, 32(3), 468-482.
- [39] White, H. (2014). *Asymptotic theory for econometricians*. Academic press.
- [40] Wooldridge, J. M. (2003). Cluster-sample methods in applied econometrics. *American Economic Review*, 93(2), 133-138.
- [41] Wooldridge, J. M. (2005). Fixed-effects and related estimators for correlated random-coefficient and treatment-effect panel data models. *Review of Economics and Statistics*, 87(2), 385-390.
- [42] Zellner, A. (1962). An efficient method of estimating seemingly unrelated regressions and tests for aggregation bias. *Journal of the American Statistical Association*, 57(298), 348-368.
- [43] Zellner, A. (1971). *An introduction to Bayesian inference in econometrics*. Wiley, New York.