



Bayesian small area estimation for skewed business survey variables

Journal:	<i>Journal of the Royal Statistical Society</i>
Manuscript ID	JRSS-OA-SC-May-17-0010.R1
Manuscript Type:	Original Article - Series C
Date Submitted by the Author:	n/a
Complete List of Authors:	Fabrizi, Enrico; Università Cattolica del S. Cuore, DISES ferrante, Maria; University of Bologna, of Statistical Sciences Trvisano, Carlo; University of Bologna,
Keywords:	Log-normal distribution, Variance-Gamma distribution, robust estimation, local shrinkage priors, regional studies

SCHOLARONE™
Manuscripts

Only

Bayesian small area estimation for skewed business survey variables

Enrico Fabrizi

Università Cattolica del S. Cuore, Piacenza, Italy – enrico.fabrizi@unicatt.it

Maria Rosaria Ferrante, Carlo Trivisano

Università di Bologna, Bologna, Italy – maria.ferrante@unibo.it, carlo.trivisano@unibo.it

Abstract

Regional and sectoral economic decision and policy rely on accurate business statistics on sub-national regions and business classes. Unfortunately, estimates based on business sample surveys can be imprecise due to the small sample sizes available for sub-populations. When this is the case, small area estimation methods can be helpful. We propose a small area technique for the estimation of totals for skewed target variables, which are typical of business data. We adopt a Bayesian approach to inference. We use a prior distribution for the random effects based on the idea of local shrinkage, which is particularly suitable when auxiliary variables with strong predictive power are available, another feature often displayed by business survey data. The proposed method is easily implemented using MCMC software. We discuss an application based on data from the Italian survey on Small and Medium Enterprises, where we estimate the total value added for subsets of a firm population obtained by cross-classifying industry, region and firm size. A simulation exercise explores the frequentist properties of the proposed estimator.

Keywords: robust estimation, Log-Normal distribution, local shrinkage priors, regional studies, Variance Gamma distribution.

1. Introduction

Regional economic decisions and policies rely on accurate business information regarding sub-national regions and business categories. The relevance of regional estimates of business aggregates and the interest in regional disparities in terms firm competitiveness and productivity is

1 demonstrated by the growing number of scientific articles in this field (see Breinlich et al., 2014 for
2 a review).

3 Regional statistics are produced by the National Statistical Institutes, and governments use
4 them to coherently allocate funds (for examples of this, see OECD, 2013, Eurostat, 2011, Eurostat,
5 2015a). For instance, the (gross) value added, that is, the total value of new goods produced and
6 services provided in a given time period, is routinely estimated at the national and sub-national
7 levels. For the EU, Eurostat releases regional estimates of the value added at levels as detailed as
8 the EU NUTS 3 regions (following the Nomenclature of Territorial Units for Statistics, Eurostat,
9 2015b) and industries (NACE Rev. 2, 1 digit, following the Statistical classification of economic
10 activities in the European Community). Sub-national estimates of value added would be even more
11 informative if they were disaggregated both in terms of industry and firm size for the purpose of
12 measuring the relative contribution of an industry and of certain firm-size classes to the regional
13 economy. Unfortunately, sample sizes of official business surveys are too small for the standard
14 design-based estimators (known as “direct estimators”) to be sufficiently precise in small domains.

15 This limitation can be overcome by model-based small area estimation methods. The small
16 area estimation literature has until very recently focused largely on the analysis of social surveys,
17 with estimation goals such as the poverty mapping (see Pfeffermann, 2014 and Pratesi, 2016 for a
18 review) and few applications for business statistics. In the last few years, awareness of this field of
19 application has grown (Burgard et al. 2014; Ferrante and Trivisano, 2010; Militino et al., 2015), as
20 well as the availability of reliable administrative archives for firms that can be used to obtain
21 auxiliary information.

22 Small area models may be broadly classified into *area level* and *unit level*. In area level
23 models, survey weighted (direct) estimates obtained for each domain are related with auxiliary
24 information at the same level of population disaggregation. In unit level models, the target variables
25 and auxiliary variables are related at the statistical unit level. Area level models straightforwardly
26 incorporate information on the sampling design and non-response re-weighting adjustments, leading
27 to design-consistent estimators whenever direct estimators are design-consistent (Rao, 2003, p.
28 117). Design consistency is a general purpose form of protection against model failures, as it
29 guarantees that, at least for large domains, estimates make sense even if the assumed model
30 completely fails. Area level modelling is less demanding in terms of data disclosure and overcomes
31 potential problems of record linkage between the survey sample and the administrative archive. For
32 these reasons, area level models will be considered in this paper.

33 Many business survey variables are positive and positively skewed (Rivière, 2002), so
34 normality is not a tenable assumption in most of the cases. Log transformation can then be

1
2
3 1 introduced in order to apply normal linear mixed models on the log scale. Predictions on the
4
5 2 original data scale require back-transformation that is a potential source of bias. Positive skewness
6
7 3 of survey variables may cause estimators of means and totals to have non-normal (positively
8
9 4 skewed) sampling distributions, when calculated on small samples. Literature on area level
10
11 5 modelling on the log scale include Fay and Herriot (1979) and Slud and Maiti (2006) that both
12
13 6 consider an empirical Bayes approach to inference.

14
15 7 In this paper we propose a full Bayes approach, accounting for all sources of uncertainty,
16
17 8 effectively dealing with back-transformation bias and implementable with widely available MCMC
18
19 9 software.

20
21 10 When predicting means or totals for business survey variables, strong covariates from
22
23 11 administrative archives are often available. For instance, in our application, aimed at predicting
24
25 12 gross value added at the domain level, we can exploit the knowledge of turnover for each firm in
26
27 13 the population. Area level totals of turnover are strongly correlated with those of value added.
28
29 14 Nonetheless, a minority of the areas will typically deviate from the relationship that characterize
30
31 15 most of the others. If we think of modelling in terms of mixed models, we have that random effects
32
33 16 would be needed for a subset of the areas (Datta, 2011) or alternatively that there are subsets of
34
35 17 random effects characterized by different variances. The specification of spike-and-slab priors can
36
37 18 be useful in this case (Datta and Mandal, 2015).

38
39 19 This paper contributes to the small area literature by proposing an approach based on local
40
41 20 shrinkage priors for the random effects (Frühwirth-Schnatter and Wagner, 2010). where spike-and-
42
43 21 slab priors are replaced by continuous gamma scale mixture of normal distributions (Griffin and
44
45 22 Brown, 2010) that lead to marginal variance-gamma distributions for the random effects. This
46
47 23 flexible modelling of random effects lead to predictions close to those we can obtain using standard
48
49 24 priors for non-outlying areas, and to less biased predictors for the areas that can be labelled as
50
51 25 outliers.

52
53 26 The strategy we propose may be applied to estimating business totals based on any
54
55 27 positively skewed variables: value added, turnover, labor cost, and income from sales and services
56
57 28 as well as the components of these main aggregates. We discuss the proposed model with reference
58
59 29 to real survey data and, more specifically, to the estimation of the total value added, giving
60
30 consideration to the fact that the value added is the input for calculating important economic
31
32 31 aggregates and performance indicators. We address the sub-populations of Italian small- and
33
34 32 medium-sized manufacturing firms classified according to sub-national region, industry and firm-
35
36 33 size classes. We limit our attention to small and medium enterprises, that is, on firms with less than
37
38 34 100 employees because in general, as well as in Italy, larger firms are censused, and small area

1 estimation is therefore not needed. We use data on the Small and Medium Enterprises sample
 2 survey (1-99 employees) conducted by the Italian National Statistical Institute (ISTAT), which
 3 provided us with this information within the framework of the BLUE-ETS project; this project has
 4 been financially supported by the EU Commission within the 7th Framework Programme.

5 The paper is organized as follows. Model specification is described in section 2. Specifically
 6 in section 2.1 closed formulas for posterior means conditionally on variance components are
 7 illustrated as posterior means are proposed as point predictors. In section 3, we apply our
 8 methodology on real survey data,. We motivate the recourse to log normal likelihood for the direct
 9 estimators. In section 4, we introduce a simulation exercise to explore the frequentist properties of
 10 the proposed predictor in comparison with some alternatives, including the estimator of Slud and
 11 Maiti (2006). Section 5 presents the study's conclusions.

12 13 14 15 **2. Small Area Estimation Model**

16 Let Y be the target variable, which we assume positive with a positively skewed distribution.
 17 Assume that Y is defined on a population U of N units, partitioned into a set of m non-overlapping
 18 domains of size N_d ($d = 1, \dots, m; N = \sum_{d=1}^m N_d$). A random sample of overall size n is taken using a
 19 possibly complex design: samples of sizes n_d are drawn from each domain. The small area nature
 20 of the problem lies in n_d being too small to allow for reliable inference for most of the domains.
 21 We assume that individual weights w_{dj} $j = 1, \dots, n_d$ are attached to responses y_{dj} to account for
 22 unequal selection probabilities and possibly other selection adjustments.

23 The normal distribution is not suitable to describe either the distribution of Y in the
 24 population nor the sampling distribution of the domain totals' direct estimators $\hat{Y}_d = \sum_{j=1}^{n_d} w_{dj} y_{dj}$.
 25 Although these are linear combinations of individual observations and can be assumed to be
 26 approximately normally distributed in large samples, in samples of small size, the sum of a few
 27 positively skewed variables remains positively skewed. We assume that the total direct estimators
 28 are log-normally distributed:

$$29 \hat{Y}_d | \theta_d, V_d \sim LN([\theta_d], [V_d]) \quad [1]$$

30 where $[\cdot]$ is used to denote a parametrization in terms of mean and variance of the distribution.
 31 Exact or approximate design-unbiasedness of totals' estimators is typical in survey sampling. The
 32 distributional assumption in [1] can be motivated directly assuming the log-normality of Y . Log-

1 normal approximations of sums of independent log-normals are justified by several authors (e.g.,
 2 Fenton, 1960; Cobb et al., 2012). Moreover, Mazmanyán et al. (2009) proposed a log-normal
 3 central limit theorem for the approximation of the sum of positively skewed random variables,
 4 although not necessarily log-normal. Eventually, the assumption of normality on the log scale when
 5 dealing with mean or total estimators of skewed variables is common in the small area literature (as
 6 in Fay and Herriot, 1979).

7 On the log-scale, a specification consistent with the *sampling model* [1] is given by:

$$8 \log(\hat{Y}_d) | \eta_d, \delta_d \sim N(\eta_d - \delta_d/2, \delta_d) \quad [2]$$

9 where $\eta_d = \ln(\theta_d)$ and $\delta_d = \text{Var}\{\log(\hat{Y}_d)\}$. $E\{\log(\hat{Y}_d)\} = \eta_d - \delta_d/2$ is in line with assuming the
 10 availability of an unbiased estimator on the original scale of the data: if $E(\hat{Y}_d) = \theta_d$, then
 11 $E\{\log(\hat{Y}_d)\} < \log(\theta_d)$. Note that $V_d = \exp\{2\theta_d + \delta_d\}[\exp\{\delta_d\} - 1]$ will depend on both parameters
 12 of the lognormal distribution.

13 In the small area literature, variances associated with direct estimators are usually treated as
 14 known constants. In practice, estimates obtained with methods such as linearization or bootstrap are
 15 smoothed using a model involving unknown parameters. In line with the literature on area-level
 16 models, we will assume that variances on the log-scale are known and denote them as $\hat{\delta}_d$.

17 We assume a multiplicative *linking model* for θ_d that links the outcome parameter to the
 18 auxiliary information in order to improve the direct estimators:

$$19 \theta_d = \exp(\mathbf{x}'_d \boldsymbol{\beta} + u_d) \quad [3]$$

20 The p -row vector \mathbf{x}'_d contains the covariates known for domain d from external sources, and u_d is a
 21 random intercept associated to θ_d . Let us assume that $u_d \sim N(0, \psi_d)$, which implies

$$22 \theta_d \sim LN(\mathbf{x}'_d \boldsymbol{\beta}, \psi_d) \quad [4]$$

23 or, equivalently, $\eta_d \sim N(\mathbf{x}'_d \boldsymbol{\beta}, \psi_d)$. We denote the model defined by *sampling model* [1] and *linking*
 24 [4] as the LN-LN model.

25 26 27 **2.1 Analysis conditional on the variance components**

28 To analyze the model defined by [1] and [4], note first that, assuming $\hat{\delta}_d$ as known
 29 ($\delta_d = \hat{\delta}_d$) we can re-write [2] as $\hat{Z}_d \sim N(\eta_d, \hat{\delta}_d)$, where $\hat{Z}_d = \log \hat{Y}_d + \frac{1}{2} \hat{\delta}_d$. We can use standard

1 results from the analysis of linear mixed models (see Rao 2003, chapter 5) to prove that,
2 conditionally on the regression coefficients $\boldsymbol{\beta}$ and the variances $\boldsymbol{\psi}_d$:

$$3 \quad \eta_d \mid \boldsymbol{\beta}, \boldsymbol{\psi}_d, data \sim N(\hat{\boldsymbol{\eta}}_d^{B1}, g_{1,d})$$

4 where $\hat{\boldsymbol{\eta}}_d^{B1} = \mathbf{x}'_d \boldsymbol{\beta} + \gamma_d (\hat{Z}_d - \mathbf{x}'_d \boldsymbol{\beta})$, $g_{1,d} = \gamma_d \hat{\delta}_d$ and $\gamma_d = \frac{\boldsymbol{\psi}_d}{\boldsymbol{\psi}_d + \hat{\delta}_d}$. Note that, as a function of the

5 shifted direct estimates \hat{Z}_d , $\hat{\boldsymbol{\eta}}_d^{B1}$ is a convex linear combination of a direct (\hat{Z}_d) and a synthetic
6 component ($\mathbf{x}'_d \boldsymbol{\beta}$), known as the linear composite estimator in the small area literature. If we
7 assume quadratic loss and define $\hat{\boldsymbol{\theta}}_d^{B1} = E(\boldsymbol{\theta}_d \mid \boldsymbol{\beta}, \boldsymbol{\psi}_d, data)$ as the point predictor for $\boldsymbol{\theta}_d$, we have
8 that

$$9 \quad \begin{aligned} \hat{\boldsymbol{\theta}}_d^{B1} &= \exp\left\{\mathbf{x}'_d \boldsymbol{\beta} + \gamma_d (\hat{Z}_d - \mathbf{x}'_d \boldsymbol{\beta}) + \frac{1}{2} \gamma_d \hat{\delta}_d\right\} \\ &= \exp\left\{\mathbf{x}'_d \boldsymbol{\beta} + \gamma_d (\log \hat{Y}_d - \mathbf{x}'_d \boldsymbol{\beta}) + \gamma_d \hat{\delta}_d\right\} \end{aligned} \quad [5]$$

10 This predictor is the product between $\exp(\hat{\boldsymbol{\eta}}_d^{B1})$ and a factor that corrects for the main bias term in
11 the back-transformation; it is in line with formula (4) of Slud and Maiti (2006).

12 It can also be shown that

$$13 \quad E(\hat{\boldsymbol{\theta}}_d^{B1}) = \exp\left\{\mathbf{x}'_d \boldsymbol{\beta} + \frac{1}{2} \boldsymbol{\psi}_d\right\} = E_M(\boldsymbol{\theta}_d) \quad [6]$$

14 where $E(\hat{\boldsymbol{\theta}}_d^{B1})$ is the expectation taken with respect to both linking and sampling models, while
15 with $E_M(\cdot)$, we denote the expectation with respect to linking model [4]. This latter result means
16 that $\hat{\boldsymbol{\theta}}_d^{B1}$ is an unbiased predictor of $\boldsymbol{\theta}_d$ in the same sense that BLUP predictors are unbiased: the
17 unconditional frequentist expectation of the estimator and the expectation of the estimand under the
18 linking model are the same. A proof of [6] can be found in Appendix 1.

19 If we remove the conditioning on $\boldsymbol{\beta}$ and assume a non-informative flat prior on $\boldsymbol{\beta}$, i.e.,

20 $p(\boldsymbol{\beta}) \propto \mathbf{1}$, then we have that

$$21 \quad \eta_d \mid \boldsymbol{\psi}_d, data \sim N(\hat{\boldsymbol{\eta}}_d^{B2}, g_{1,d} + g_{2,d})$$

22 where $\hat{\boldsymbol{\eta}}_d^{B2} = \mathbf{x}'_d \hat{\boldsymbol{\beta}}_{gls} + \gamma_d (\hat{Z}_d - \mathbf{x}'_d \hat{\boldsymbol{\beta}}_{gls})$, $\hat{\boldsymbol{\beta}}_{gls} = \left(\sum_d \frac{1}{\boldsymbol{\psi}_d + \hat{\delta}_d} \mathbf{x}_d \mathbf{x}'_d\right)^{-1} \sum_d \frac{1}{\boldsymbol{\psi}_d + \hat{\delta}_d} \mathbf{x}'_d \log \hat{Y}_d$,

23 $g_{2,d} = (1 - \gamma_d)^2 \mathbf{x}'_d \left(\sum_d \frac{1}{\boldsymbol{\psi}_d + \hat{\delta}_d} \mathbf{x}_d \mathbf{x}'_d\right)^{-1} \mathbf{x}_d$. As a consequence, the point predictor under quadratic

24 loss will be given by

$$\hat{\theta}_d^{B2} = \exp \left\{ \mathbf{x}'_d \hat{\boldsymbol{\beta}}_{gls} + \gamma_d \left(\hat{Z}_d - \mathbf{x}'_d \hat{\boldsymbol{\beta}}_{gls} \right) + \frac{1}{2} (g_{1,d} + g_{2,d}) \right\} \quad [7]$$

(see Appendix 1 for a proof). Unlike the empirical Bayes approach advocated by Slud and Maiti (2006), who plug estimates of unknown parameters into [7], a full Bayes approach accounts for the effect that the extra-variation implied by the estimation of $\boldsymbol{\beta}$ has on the point predictor; in fact, the expectation of a log-normal variable depends on both the expectation and variance on the log scale.

To fully account for all sources of uncertainty, we should remove the conditioning on the variance components ψ_d ; unfortunately, for sensible choices of the prior, this leads to posterior distributions for θ_d that cannot be written in closed form and should therefore be explored by means of computational algorithms such as the Markov Chain Monte Carlo considered in this paper.

2.2 The distribution for the random effects and specification of hyperpriors

The main difference between [4] and the linking model adopted by most of the small area literature on Fay-Herriot type models (Jiang and Lahiri, 2005; Pfeiffermann, 2014) is that the variances associated with random intercepts are in [4] domain-specific, implying local shrinkage instead of the ordinary global shrinkage that we would have had assuming $\psi_d = \psi$, $\forall d$. In a different context, the specification of a distribution for random intercept based on local shrinkage is discussed in Frühwirth-Schnatter and Wagner (2010).

Datta et al. (2011) note that in the presence of good covariates, the variability of the small area parameters may be accounted for by a synthetic estimator, and the inclusion of a random effect term may be unnecessary. When random effects are needed for a subset of the areas, the specification of spike-and-slab priors can be useful (see Datta and Mandal, 2015). Spike-and-slab priors amount to assuming that random intercepts are sampled from a mixture of two normal distributions.

When analyzing business data, it is quite likely that auxiliary variables with strong predictive power are available. When this is the case, the bulk of the direct estimates will be well fitted by the synthetic model (without random intercepts), so the associated ψ_d are likely to be small, with a minority of areas that require larger area-specific intercepts (and thereby larger ψ_d).

Our specification for the distribution of u_d , $d = 1, \dots, m$ is based on infinite mixtures of normal distributions. Following Griffin and Brown (2010), our specification uses *Gamma* mixtures of normal distributions. Specifically, we assume:

$$u_d | \psi_d \stackrel{ind}{\sim} N(0, \psi_d) \quad [8]$$

$$\psi_d | a, \lambda \stackrel{ind}{\sim} \text{Gamma}(a, \lambda) \quad [9]$$

$$\lambda | b_0, c_0 \sim \text{Gamma}(b_0, c_0) \quad [10]$$

This leads to spiked priors for the random effects that at the same time have tails that are heavier than those of the normal distribution. Griffin and Brown (2010) observe that for small values of the shape parameter a , the prior specification [8]-[10] leads to a marginal prior distribution for u_d that mimics the behavior of spike-and-slab priors based on finite mixtures. This infinite mixture specification is computationally easier to deal with.

Other choices for the mixing distribution such as the popular *Inverse Gamma* would lead to platikurtic distributions with heavy tails, such as those in the t family; this contrasts with the idea of severe shrinkage for most of the areas, which is consistent with a large probability mass close to 0.

Specifically, prior specification [8]-[10] implies that $u_d | a, \lambda$ follows a Variance Gamma distribution, i.e.,

$$u_d \sim VG(a, \sqrt{2\lambda}, 0, 0)$$

(see Bibby and Sørensen, 2003 for more details on this distribution). This marginal prior distribution is symmetric and has expectation $E(u_d) = 0$ and variance $V(u_d) = a / \lambda$. It belongs to the family of generalized hyperbolic distribution (Barndorff-Nielsen, 1977). The conjugate hierarchy in [8]-[10] also facilitates MCMC sampling.

In line with Frühwirth-Schnatter and Wagner (2010), we set the shape parameter a to a fixed (small) value, while we specify a prior on the global parameter λ . As far as a is concerned, we focus on two choices, $a = 1$, $a = 0.5$.

The choice $a = 1$ implies that u_d is a priori distributed as a double-exponential or Laplace, which, combined with the normal prior conditional on ψ_d , recalls the Bayesian lasso of Park and Casella (2008).

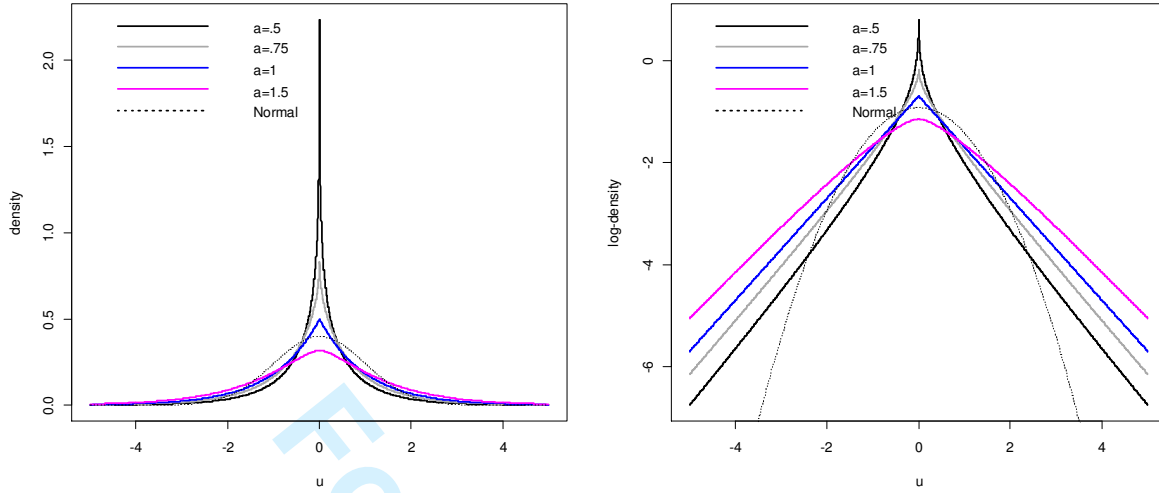


Figure 1: density (left) and log-density (right) functions of the Variance Gamma distribution.

$VG(a, \sqrt{2\lambda}, 0, 0)$ for $\lambda=1$ and various values of a .

The case $a = 0.5$ represents a more peaked prior distribution and encourages more shrinkage towards 0 of small random intercepts (see Figure 1). Moreover, it leads to a *half-t* marginal prior on $\sqrt{\psi_d}$. The *half-t* prior for standard deviations is discussed in Gelman (2006) and recommended whenever it makes sense to put a sizeable mass of prior probability close to 0. It can be shown that once $\boldsymbol{\psi} = \{\psi_d\}$ and $\boldsymbol{\tau} = \{\sqrt{\psi_d}\}$, $d = 1, \dots, m$, are defined, under prior [8]-[10] and $a = 0.5$,

$$\boldsymbol{\tau} | b_0, c_0 \sim Mht\left(0, \frac{2b_0}{c_0} \mathbf{I}, 2b_0\right). \quad [11]$$

With $Mht(\cdot)$ we denote the multivariate distribution (with support \mathbb{R}^{+n}) that is obtained from a multivariate t distribution applying the absolute value transformation on each component of the random variable. We can also prove that each prior $\sqrt{\psi_d} | b_0, c_0$ is a (univariate) *half-t* and the priors for two different variance components are uncorrelated, i.e. $Cov(\psi_d, \psi_{d^*} | b_0, c_0) = 0$ whenever $d \neq d^*$. See appendix 1 for a proof of [11] and the other statements.

As for the prior specification of the remaining parameters, diffuse independent normal priors can be specified for the components of $\boldsymbol{\beta}$. We can set $b_0 = 2$, which implies that $E(\lambda^{-1}) = c_0$. This helps to interpret c_0 as a scaling constant for the random effects variance $V(u_d) = a / \lambda$. The choice of c_0 depends on the scale for the random effects in the problem being considered. According to

[11], the parameter b_0 can be interpreted in terms of degrees of freedom of the marginal prior $p(\boldsymbol{\tau})$, so the choice of $b_0 = 2$ is in line with selecting half- t priors with a very small number of degrees of freedom (Gelman, 2006).

3. An Application for Italian Small and Medium Enterprise Survey Data

In this section, we illustrate the methodology we discussed using real survey data. We use data on the Small and Medium Enterprises (SME) sample survey, wave 2008, conducted by the Italian National Statistical Institute (ISTAT). Specifically, we target the estimation of the total value added (VA) for small domains of the population of Italian small and medium manufacturing firms (less than 100 employees). The domains we focus on are smaller than those for which ISTAT provides reliable estimates. Specifically, our domains are defined by cross-classifying: the 20 Italian NUTS 2 administrative regions, the economic industrial sector (NACE Rev. 2, 2 digit, 22 industries) and firm size (4 classes: less than 10 employees, 10 to 19 employees, 20 to 49 employees, 50 to 99 employees). As anticipated, for domains as small as those that we target, standard design-based estimators are characterized by unacceptably large variances.

3.1 Direct estimators and sampling model

The SME survey sampling design is stratified and strata are defined by cross-classifying NACE Rev. 2 (4 digits) Italian administrative regions and company size in the four classes defined above. The domains we are interested in are planned because they are unions of sampling strata.

Let \hat{Y}_{ijr} be the direct estimator of the parameter θ_{ijr} , where i indexes the economic activity ($i = 1, \dots, 22$), j the size classes ($j = 1, \dots, 4$) and r the regions ($r = 1, \dots, 20$). The potential number of 1760 domains falls to 1165, as some of the populations obtained by cross-classification are empty and some very small. We excluded all the domains characterized by a sampling rate over 0.75.

The actual sample sizes for the domains we consider ranges from 2 to 184, with a median of 8, a mean of 13.5 and .75 and .9 quantiles equal to 16 and 30 respectively.

Direct estimates can be obtained using the calibration estimator that ISTAT adopts for the SME survey. Calibration estimators can be written as weighted sums. ISTAT's published weights are obtained by multiplying base weights (the inverse of the inclusion probabilities) by factors adjusting for non-response and calibrating to known totals. Let the estimated total be denoted as

$$\hat{Y}_{ijr} = \sum_{k \in d_{ijr}} w_{ijr,k} y_{ijr,k}, \text{ where } y_{ijr,k} \text{ is the value added of the } k\text{-th firm in sector } i, \text{ size class } j, \text{ region } r.$$

1 We assume that $E(\hat{Y}_{ijr}) = \theta_{ijr}$ with $Var(\hat{Y}_{ijr}) = V_{ijr}$. We estimate V_{ijr} using linearization-based
 2 variance estimators and denote these estimates as \hat{V}_{ijr} .

3 In our sampling model we assume log-normality according to [1]. To justify this assumption
 4 for our data we proceed in two steps: first, we check whether log-normality is a sensible assumption
 5 for domain specific sample data, then we use a simple simulation exercise to assess whether log-
 6 normality is to be preferred to normality as the sampling distribution of total estimators given the
 7 sample sizes we have in our analysis.

8 For all the domains with $n_{ijr} \geq 3$, we tested normality and log-normality using the Shapiro-
 9 Wilk test. Results are reported in table 1 below. In reporting the results we consider separately the
 10 smallest 90% of the domains $n_{ijr} \leq 30$ and the largest 10%. In the smaller domains, for which the
 11 test is relatively less powerful, both normality and log-normality tend to be not rejected, but
 12 normality fails clearly more often. In larger domains, when the test has more power, normality is
 13 rejected in the large majority of cases, while log-normality is accepted in more than 70% of the
 14 cases.

15
 16 Table 1: *Checking normality and log-normality within domain-specific samples using Shapiro-Wilk
 17 test. Percentage of non-rejections at the 0.01 significance level*

n_{ijr}	normality	log-normality
≤ 30	0.733	0.959
> 30	0.087	0.713
overall	0.672	0.943

18
 19 From table 1 we conclude that log-normality is a sensible assumption for the distribution of
 20 VA within domains. We actually assume that direct estimators are log-normally distributed
 21 according to the arguments illustrated in previous section. To check this, we consider a set of log-
 22 normal populations: $Y_i \sim LN(\tilde{\mu}_i, \tilde{\sigma}_i)$, $i = 1, \dots, L$ where $L = 77$ is the number of domains with
 23 $n_{ijr} > 30$ for which log-normality was not rejected and $\tilde{\mu}_i, \tilde{\sigma}_i$ are the parameters according to
 24 maximum likelihood for these domains. For each of these populations we generated simple random
 25 samples of sizes $R = 10,000$ for each of the following sample sizes: $n_i = 5, 10, 15, 20$. Note that 20
 26 represent the 0.8 quantile of the sample size distribution in our application.

27 We evaluate how far is the empirical sampling distribution of the sample mean from the
 28 normal and the log-normal distributions in terms of Kolmogorov-Smirnov distance averaging over

the L populations. In fact, formal hypothesis testing of distributional assumptions with a sample of replicates as large as 10,000 would lead to rejections even in presence of negligible departures from the null. Results, summarized in table 2, show how log-normality is to be preferred to normality for all sample sizes. We can also note that as the sample size grows larger, the difference between the two distances decreases.

Table 2: Kolmogorov-Smirnov distances between the Monte Carlo ($R = 10,000$) distribution of the sample mean and two reference distributions, for various sample sizes.

Reference distribution	sample size (n_i)			
	5	10	15	20
Log-normal	0.012	0.014	0.015	0.015
Normal	0.119	0.097	0.081	0.073

In order to obtain more stable direct variance estimates, we smooth them through the Generalized Variance Function approach (Wolter, 1986). To begin with, we consider that under the log-normality assumption introduced in [1], we have that

$$\text{Var}\{\log(\hat{Y}_{ijr})\} = \log\{CV^2(\hat{Y}_{ijr}) + 1\}. \quad [12]$$

Thus, the smoothing can be conducted on $CV^2(\hat{Y}_{ijr}) = \frac{\hat{V}_{ijr}}{\hat{Y}_{ijr}^2}$. After careful explorative analysis, we assume that $CV^2(\hat{Y}_{ijr})$ varies with the size class (j) but not with economic activity (i) or with regions (r). This leads to the following smoothing equation for the direct estimate of V_{ijr} , \hat{V}_{ijr} :

$$\hat{V}_{ijr} = \phi_j \frac{\hat{Y}_{ijr}^2}{n_{ijr}} \left(1 - \frac{n_{ijr}}{N_{ijr}}\right) + v_{ijr} \quad [13]$$

with $E(v_{ijr}) = 0$, $V(v_{ijr}) = \kappa$ and where a finite population correction factor is also considered to account for varying and occasionally non-negligible sample rates. The parameter ϕ_j can be interpreted as the smoothed squared coefficient of variation multiplied for the size of the domain n_{ijr} . The domain sample size n_{ijr} in the denominator of [13] allows for the decrease in the coefficient of variation when the sample size increases. Smoothed squared estimated coefficient of

variations are given by $CV_{smooth}^2(\hat{Y}_{ijr,k}) = \frac{\phi_j}{n_{ijr}} \left(1 - \frac{n_{ijr}}{N_{ijr}}\right)$; the first, second and third quartiles of

1 $CV^2_{smooth}(\hat{Y}_{ijh,k})$ estimated on our data set are 31%, 45% and 65%, respectively. These results
 2 confirm the need to adopt a small area model approach.

3 We can then adapt the sampling model [1] to our problem as

$$4 \log \hat{Y}_{ijr} | \eta_{ijr}, \hat{\delta}_{ijr} \sim N\left(\eta_{ijr} - \frac{1}{2} \hat{\delta}_{ijr}, \hat{\delta}_{ijr}\right) \quad [14]$$

5 where $\eta_{ijr} = \log(\theta_{ijr})$ and $\hat{\delta}_{ijr} = \log\{CV^2_{smooth}(\hat{Y}_{ijr}) + 1\}$ according to [12].

8 3.2 Auxiliary information and linking model

9 As an auxiliary variable, the log total turnover in each domain is available. This auxiliary
 10 information refers to the Italian firms' population and it is provided by the Italian Statistical
 11 Register of Active Enterprises-ASIA archive. The predictive power of this covariate is quite strong:
 12 the squared correlation coefficient is equal to 0.87 when calculated on variables on their original
 13 scale, and it is equal to 0.79 for the log transformations. In the original scale the high correlation
 14 level is influenced by few observations with a larger scale with respect to most of the others.

15 We assume the multiplicative *linking model* [4] for θ_{ijr} to link the outcome parameter to the
 16 auxiliary information

$$17 \eta_{ijr} \sim N(\beta_0 + x_{ijr}\beta_1, \psi_{ijr}) \quad [15]$$

18 where x_{ijr} is the log-total turnover for the domain in question. Equivalently, we can write
 19 $\eta_{ijr} = \beta_0 + \text{ltt}_{ijr}\beta_1 + u_{ijr}$. The prior for the vector of domain-specific random intercepts u_{ijr} is
 20 specified according to [8]-[10]. As for the prior specifications not already discussed, we set
 21 $\beta_0 \sim N(0, 10^5)$, $\beta_1 \sim N(0, 10^5)$, $b_0 = 2$, $c_0 = 1$. We chose these values as they provide a reasonable
 22 scale for the random effects variance in our problem.

23 We also consider the LN-LN model with an alternative choice for the prior distribution on
 24 u_{ijr} :

$$25 u_{ijr} | \sigma^2 \sim N(0, \sigma^2), \sigma^2 \sim \text{InverseGamma}(c, d) \quad [16]$$

26 This prior specification, which implements global shrinkage, can be considered as a
 27 benchmark for evaluating the effects of prior specification approximating spike-and-slab introduced
 28 in the previous section, and it represents a routine choice in many applications. We set $c = 0.01$,
 29 $d = 0.01$.

3.3 Markov Chain Monte Carlo computational issues

Parameter estimates are obtained by summarizing the posterior distributions approximated by the output of Markov Chain Monte Carlo (MCMC) integration via the Gibbs sampling algorithm. By assuming a quadratic loss, the posterior means are adopted as estimates of the area specific parameters. Posterior variances are used as a measure of uncertainty. To carefully assess the convergence, we run three parallel chains of 25,000 runs each, the starting point being drawn from an over-dispersed distribution. The convergence of the Gibbs sampler was monitored by visual inspection of the chains' plots and autocorrelation diagrams and by means of the potential scale reduction known as the Gelman-Rubin statistic (Carlin and Louis, 2000, ch. 5). Both models displayed fast convergence; we discarded the first 5,000 iterations from each chain. To obtain estimates, we used the OpenBugs software package, which can be downloaded for free on the internet and it is open source.

3.4 Comparing alternative models

In order to choose among competing models, we compute the Deviance Information Criterion (DIC) and the logarithm of the pseudo-marginal likelihood (LPLM, Ibrahim et al., 2001). A model is preferred if it displays a lower DIC value. Table 3 reports the DIC results for the whole set of small area models estimated. DIC values show that, in line with expectations, the log-normality assumption at the sampling level performs better in terms of DIC with respect to the model assuming normality. The ordering of alternative models using LPLM is consistent with that obtained using the DIC. The adoption of the Variance Gamma for the random intercepts u_{ijr} leads to a further reduction in DIC with respect to the more common specification [16].

Table 3: Comparison of alternative assumptions on the distributions of the random effects

Shrinkage	Prior on random intercepts u_{ijr}	a	DIC	LPLM	median CVR
global	[16]	-	15340	-7846	0.391
local	[10]-[12]	1	15230	-7808	0.421
local	[10]-[12]	0.5	15220	-7798	0.455

1
2
3 1 We also compare the median reduction of the coefficient of variation of estimators with
4
5 2 respect to the direct ones, defined as $median(CVR_k^h)$. CVR_k^h is defined as $CVR_k^h = 1 - CV_k^h / CV_k^{DIR}$,
6
7 3 where the CV_k^h is the coefficient of variation calculated on the posterior of θ_k (k being a generic
8
9 4 index for the areas) under model h , while CV_k^{DIR} is the coefficient of variation of the direct
10
11 5 estimators calculated from the randomization distribution.

12
13 6 The posterior predictive approach can be used to assess the fit of a model (Gelman et al.,
14
15 7 1996). We consider a discrepancy measure suggested in the context of small area estimation by You
16
17 8 and Rao (2002) and considered also in Fabrizi and Trivisano (2016):

$$18 \quad dis_{ijr} = P(\hat{Y}_{ijr} < Y_{ijr}^*)$$

19
20 10 where Y_{ijr}^* is generated from the posterior predictive distribution. The discrepancy measure is aimed
21
22 11 at checking whether, for each area, the posterior predictive distribution is centered around the direct
23
24 12 estimate. Values of dis_{ijr} far from 0 and 1 would provide evidence of systematic under or over
25
26 13 estimation. For the log-normal model endowed with priors [10]-[12] and $a = 0.5$ (i.e. the best
27
28 14 model in terms of DIC and LPLM), we have that the average of the discrepancy measure over the
29
30 15 set of areas is 0.499 with 0.25 and 0.75 quantiles equal 0.32 and 0.68 respectively, which means an
31
32 16 adequate fit. Less than 10% of the areas shows dis_{ijr} out of the range (0.2, 0.8). Similar summaries
33
34 17 can be obtained for the other models considered in table 3.

35 18 Results on $median(CVR_k^h)$, reported in table 3, highlight that the whole set of considered
36
37 19 small area estimators considerably reduce the variability of direct estimators, which is consistent
38
39 20 with the availability of a strongly predictive auxiliary variable. Nonetheless, even if exploiting the
40
41 21 same auxiliary information, the models perform differently, as the prior specification has a non-
42
43 22 negligible impact. Prior specifications mimicking the spike and slab behavior allow for a further
44
45 23 gain in efficiency with respect to priors ordinarily used in this type of analysis.
46
47 24

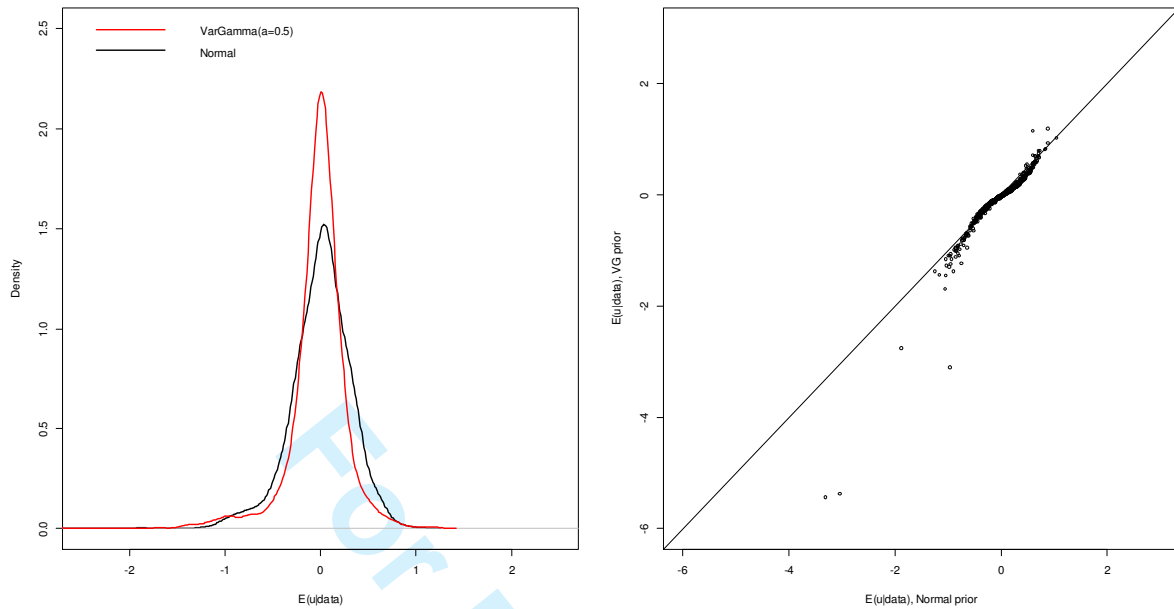


Figure 2: Left panel: kernel densities of $E(u_{ijr} | data)$ over the set of small areas under Normal and VG priors. Right panel: $E(u_{ijr} | data)$ under Normal prior vs $E(u_{ijr} | data)$ under VG prior. For VG priors, $a = 0.5$ is assumed.

To evaluate the improvements allowed by the model-based proposed predictor we can compare the number of small areas with values of the coefficient of variation CV less than 16.6%, between 16.6% and 33.3% and over 33.3% for the direct and the model-based predictor. These thresholds for CV were suggested by Statistics Canada (2007) to provide quality level guidelines for publishing small area estimates; those with a coefficient of variation less than 16.6% are considered reliable for general use. Estimates with a coefficient of variation between 16.6% and 33.3% should be accompanied by a warning to users. Estimates with coefficients of variation larger than 33.3% are deemed to be unreliable. Less than 25% of the direct estimates have associated CV below the 33.3% threshold, while for the model based ones this number grows to 70%. Although the uncertainty around the small area estimates remain sizeable and not all estimates would be publishable, the application of the proposed method endows most sub-population with a publishable estimate in spite of the small sample sizes.

Figure 2 displays the impact of alternative prior specifications on the ensemble of the random intercepts' posterior means $E(u_{ijr} | data)$ under the Normal prior [16] and Variance Gamma priors $u_{ijr} \sim VG(0.5, \sqrt{2\lambda}, 0, 0)$. Based on the left panel, it is clear that the Variance Gamma leads to

1
2
3 1 a more peaked distribution of estimated random intercepts, as predicted by theory. As tail behavior
4 is difficult to read from density estimates, in the right panel of Figure 2, we plotted the point
5 2
6 3 posterior expectation under the Normal prior vs those obtained under the Variance Gamma. The
7
8 4 peak around 0 is still apparent from the inflection of the points' cloud approximately at 0; heavier
9
10 5 tails under Variance Gamma can also be appreciated: under the Normal prior, $E(u_{ijr} | data)$ lies
11
12 6 within the interval $(-4, 2)$, while they do not under the Variance Gamma specification.

13
14 7 As the purpose of small area estimators is to complement the direct estimates obtained using
15
16 8 ordinary survey weighted methods, robustness with respect to modelling assumptions is a major
17
18 9 concern. As anticipated, the recourse to area-level model-based methods entails design consistency.
19
20 10 When the area-specific sample sizes are large (and occasionally they are) the small area estimate
21
22 11 will be close to that obtained under standard design-based methods. This offers, at least for the
23
24 12 larger domains, protection against model misspecifications; moreover, it automatically guarantees
25
26 13 that, in the case of large domains, model-based and design-based estimates are automatically in
27
28 14 agreement.

17 4. A simulation assessment of the frequentist properties of the proposed point estimators

18
19 18 In this section we introduce a simulation exercise with the aim of investigating the
20
21 19 frequentist properties of point estimators (i.e. posterior means) introduced in section 4. We study
22
23 20 bias, mean square error and frequentist coverage of posterior probability intervals.

24
25 21 The simulation is design-based and we do not assume any parametric distribution when
26
27 22 generating the data. We create a synthetic population merging the samples of the 2007 and 2008
28
29 23 SME survey discussed in previous section. We drop from the 2007 wave those firms sampled also
30
31 24 in the 2008. We obtain a population of size $N = 30451$. Domains are defined by cross-classifying
32
33 25 the population by firm size and industry sector; with respect to the data analysis of previous section
34
35 26 we collapse the regions. By reducing the number of domains we create sub-populations large
36
37 27 enough to be sampled using reasonable sampling rates. Collapsing by region has a milder impact on
38
39 28 sub-population skewness with respect to firm size or industry sector. Thereby our synthetic
40
41 29 population is divided into $m = 88$ domains whose sizes N_d ($d = 1, \dots, m$) range from 14 to 1339 with
42
43 30 an average of 346. The same target parameters (total VA) and auxiliary variable (turnover) studied
44
45 31 in section 4 are considered.

46
47 32 We keep this synthetic population as fixed and we repeatedly stratified samples with
48
49 33 proportional allocation and a 4% sampling rate. Resulting domain-specific sample sizes n_d are

1 adjusted so that $\min(n_d)=3$. The resulting average domain-specific sample sizes is 14.05, with a
 2 maximum of 54. The mean sample size is very close to that of the application.

3 The Monte Carlo exercise is based on $S = 2000$ repeated samples. Direct estimates and their
 4 variances are calculated using analytic formulas. The point estimators we compare in the simulation
 5 are:

6 *i)* the posterior mean associated to the Fay-Herriot model specified on the un-transformed scale
 7 (UFH). The untransformed Fay-Herriot model can be described as follows: $\hat{Y}_d \sim N(\theta_d, V_d)$,

8 where V_d is the variance of the direct estimator \hat{Y}_d , $\theta_d \sim N(\beta_0 + \beta_1 x_d, \sigma_v^2)$, $\sigma_v \sim Unif(0, A)$,

9 $A = 1000$, $\beta_i \sim N(0, 10^5)$ $i = 0, 1$. We consider the UFH model as it is probably the most

10 “basic” Bayesian model a practitioner would think of for analyzing these data.

11 *ii)* the predictor proposed by Slud and Maiti (2006), we denote as SM;

12 *iii)* the posterior mean obtained from the log-normal model [1]-[3] endowed with the global
 13 shrinkage prior [16]. We denote the estimator as LNGS (“Log-Normal with Global
 14 Shrinkage”)

15 *iv)* the posterior mean associated to the model [1]-[3] but endowed with local shrinkage prior
 16 described in [8]-[10]; we label the estimator as LNLS (“Log-Normal with Local Shrinkage”).

17
 18 Let's denote with est_{ds} the generic estimator calculated for the domain d in the replication s . We
 19 compare alternative estimators in terms of relative bias, relative root mean square error and
 20 frequentist coverage of probability intervals based on the posterior distribution. Specifically we
 21 consider the frequentist coverage defined by the $\alpha/2$ and $(1-\alpha/2)$ quantiles posterior distribution
 22 of the target parameter (w.r.t. the coverage probability $1-\alpha$) and set $\alpha = 0.05$. Comparison tools are
 23 defined as follows:

$$24 \quad RABIAS = \frac{1}{m} \sum_{d=1}^m \left| \frac{\frac{1}{S} \sum_{s=1}^S est_{ds} - Y_d}{Y_d} \right|,$$

$$25 \quad RRMSE = \frac{1}{m} \sum_{d=1}^m \sqrt{\frac{\frac{1}{S} \sum_{s=1}^S (est_{ds} - Y_d)^2}{Y_d^2}},$$

$$26 \quad COV95 = \frac{1}{m} \sum_{d=1}^m \frac{1}{S} \sum_{s=1}^S 1\{Y_d \in (p_{\theta_{d,ldata}}(0.025), p_{\theta_{d,ldata}}(0.975))\}$$

27

As we are interested in the frequentist coverage of Bayes estimators, $COV95$ is calculated only for i, iii, iv . Results are summarized in table 4.

Table 4: *Comparison of alternative predictors based on the Monte Carlo experiment*

est	<i>RBIAS</i>	<i>RRMSE</i>	<i>COV95</i>
UFH	0.1573	0.2026	0.708
SM	0.1263	0.1733	--
LNLS	0.1240	0.1719	0.929
LNLS	0.1113	0.1723	0.941

Results from table 4 show how the posterior means based on the log-normal model with either local or global shrinkage priors and the predictor of Slud and Maiti (2006) perform very closely in terms of mean square error. In terms of bias LNLS is better, but its variance is somewhat bigger, as we can expect from a more flexible, richly parametrized model. Actually, they are based on similar ideas and models, only the priors or the way hyper-parameters are dealt with are different, so the results are in line with expectation. We did not expect the bias to be close to 0: in small area estimation you compromise between the efficiency of a biased synthetic predictor and the unbiasedness of large variance direct estimators; to some extent estimators associated to areas with very small area-specific sample sizes are naturally biased.

The naïve Fay-Herriot model, specified on the un-transformed scale, performs worse in terms of both bias and mean square error; the frequentist coverage of the posterior intervals is well below the 0.95 nominal level. This relatively poor performances reflect the misspecification of the model, based on the assumption of normality of the direct estimators. It also assumes a linear relationship between direct estimators and the auxiliary variable on the original scale of the data (instead of a linear relationship on the log scale); we already noted that this assumption is not completely unrealistic, so misspecification of the sampling model can be held as responsible for the not completely satisfying results.

The two hierarchical models lead to close performances also in terms of frequentist coverage of posterior intervals. The advantage of using local shrinkage priors can be appreciated if we consider the performance for outlying areas, that is those characterized by a deviation from the synthetic component much larger than most of the remaining areas. We investigate performances separately for the areas characterized by the larger (on average) model residual. Results related to the “worst case” area are presented in table 5. We note that this area-specific sample is $n_d = 34$, well above the average sample size of the simulation.

Table 5: Comparison of alternative predictors for area 3, characterized by the largest model residual

est	RBIAS	RRMSE	COV95
SM	0.4763	0.4790	--
LNGS	0.4414	0.4453	0.505
LNLS	0.0144	0.2515	0.947

The LNGS and the SM predictor are based on a similar global shrinkage idea. Results in table 5 show how the common variance parameter assumed for the random effects cannot accommodate the “outlier”; the associated model based estimators are severely shrunken towards the synthetic component: this implies large bias and poor frequentist coverage of the posterior intervals. Local shrinkage prior associated to LNLS is more flexible and leads to an almost unbiased predictor and good coverage.

5. Conclusions

We introduced a Bayesian methodology that is useful for small area estimation of means and totals of variables that are positively skewed. This type of variable is often encountered in business surveys. We devote special attention to the specification of a prior distribution for the random effects; our proposal, based on the idea of local shrinkage, is well suited when auxiliary variables with strong predictive power is available, a feature often displayed in business survey data.

The proposed methodology can be easily implemented using widely available MCMC software. Openbugs codes, as well as formulas for the full conditional distributions needed for an independent implementation of the algorithm, are available upon request from the authors.

In summary, we showed that the predictor based on local shrinkage prior has overall acceptable frequentist properties, comparable to the alternatives we consider in the exercise. We introduce this prior specification to deal with situations where powerful auxiliary information is available, most of the areas are well fit by the assumed model and only a minority are outlying, characterized by larger model residuals. For these areas, local shrinkage priors can lead to estimators with reduced bias and thereby more efficient.

References

- Barndorff-Nielsen, O.E. (1977). Exponentially decreasing distributions for the logarithm of particle size, *Proceedings of the Royal Statistical Society, series A*, 353, 401-419.
- Breinlich, H., Ottaviano, G.I.P., & Temple, J.R.W. (2014). Regional Growth and Regional Decline. in: *Handbook of Economic Growth*, edition 1, vol. 2, ch. 4, 683-779 Elsevier.
- Burgard, J.P., Munnich R., & Zimmermann T. (2014). The Impact of Sampling Designs on Small Area Estimates for Business Data. *Journal of Official Statistics*, 30, 4, 749-771.
- Carlin, B.P., & Louis, T.A. (2000). *Bayes and empirical Bayes data analysis*. New York, Chapman and Hall.
- Cobb B.R., Rumi R. & Salmeron A. (2012). Approximating the distribution of a sum of log-normal random variables. In: *Proc. of the VI European Workshop on Probabilistic Graphical Models*, Granada.
- Datta, G.S. & Lahiri, P. (1995). Robust hierarchical Bayes estimation of small area characteristics in presence of covariates and outliers. *Journal of Multivariate Analysis*, 54, 2, 310-328.
- Datta, G. S., Hall, P. and Mandal, A. (2011). Model selection by testing for the presence of small-area effects in area-level data. *Journal of the American Statistical Association*, 106, 362-374.
- Datta, G.S., Mandal A. (2015) Small Area Estimation with Uncertain Random Effects, *Journal of the American Statistical Association*, in press.
- Eurostat (2011), *Key figures on European business - with a special feature on SMEs*, Eurostat Pocketbooks.
- Eurostat (2015a). *Eurostat regional yearbook 2015. Statistical books, General and regional statistics*.
- Eurostat (2015b). *Regions in the European Union - Nomenclature of territorial units for statistics - NUTS 2013/EU-28*.
- Fabrizi, E., Ferrante, M.R., Pacei, S. & Trivisano C. (2011). Hierarchical Bayes multivariate estimation of poverty rates based on increasing thresholds for small domains. *Computational Statistics & Data Analysis*, 55, 1736 – 1747.
- Fabrizi, E. & Trivisano C. (2016). Small area estimation of the Gini concentration coefficient. *Computational Statistics & Data Analysis*, 99, 223 – 234.
- Fay, R. & Herriot, R. (1979). Estimates of income for small places: an application of James-Stein procedures to census data. In: *Journal of the American Statistical Association*, 74, p. 269-277.
- Fenton, L. (1960) The Sum of Log-Normal Probability Distributions in Scatter Transmission Systems. *IRE Transactions on Communication Systems* 8(1), 57-67.

1
2
3 1 Ferrante, M.R., & Trivisano C. (2010). Small area estimation of the number of firms'
4 recruits by using multivariate models for count data. *Survey Methodology*, 36, 2, 171-180.

5
6 3 Fruhwirth-Schnatter, S., & Wagner, H. (2010). Bayesian variable selection for random
7 intercept modelling of Gaussian and non-Gaussian data. In J. Bernardo, M. Bayarri, J.O. Berger,
8 A.P. Dawid, D. Heckerman, A.F.M. Smith, M. West (eds.), *Bayesian Statistics*, 9, 165–200, Oxford
9 Univ. Press.

10
11 7 Gelman A., Meng X., Stern H. (1996) Posterior predictive assessment of model fitness via
12 realized discrepancies. *Statistica Sinica*, 6, 733-807.

13
14 9 Gelman A. (2006). Prior distributions for variance parameters in hierarchical models.
15 *Bayesian Analysis*, 1, 3, 515–533.

16
17 11 Griffin, J.E., & Brown, P.J. (2010). Inference with normal-gamma prior distributions in
18 regression problems. *Bayesian Analysis*, 5, 171- 188.

19
20 13 Horrace W.C. (2005) Some results on the multivariate truncated normal distribution, *Journal*
21 *of Multivariate Statistics*, 94, 209-221.

22
23 15 Ibrahim J., Chen M., Sinha D. (2001). *Bayesian Survival Analysis*. Springer-Verlag, New
24 York.

25
26 17 Jiang J. & Lahiri P. (2005) Mixed model prediction and small area estimation (with
27 discussion), *TEST*, 15, 1-96

28
29 19 Mazmnyan, L., Ohanyan, V., Trietsch, D. (2009) The lognormal central limit theorem for
30 positive random variables. Reproduced as Appendix in Baker K.R, Trietsch D. (2009), *Principles of*
31 *sequencing and scheduling*, John Wiley and Sons, New York.

32
33 22 Militino, A.F., Ugarte, M.D. & Goicoa T. (2015). Deriving small area estimates from
34 information technology business surveys. *Journal of the Royal Statistical Society A*, 178, 1051-
35 1067.

36
37 25 Park, Casella G. (2008). The Bayesian lasso, *Journal of the American Statistical*
38 *Association*, 103, 681–686.

39
40 27 Pratesi M. ed. (2016). *Analysis of Poverty Data by Small Area Estimation*, John Wiley, New
41 York.

42
43 29 OECD (2013), *OECD Regions at a Glance 2013*: OECD Publishing.
44 http://dx.doi.org/10.1787/reg_glance-2013-en

45
46 31 Pfeffermann D. (2014). Small Area Estimation. In *International Encyclopedia of Statistical*
47 *Science*, M. Lovric (eds.), 1346-1349, Springer-Verlag.

48
49 33 Rao J.N.K. (2003). *Small area estimation*, John Wiley and Sons. New York.

- 1
2
3 1 Rivière P. (2002): What Makes Business Statistics Special? *International Statistical Review*,
4 2 70, 1, 145-159.
5
6 3 Slud, E. V. & Maiti, T. (2006). Mean-squared error estimation in transformed Fay–Herriot
7 4 models. *Journal of the Royal Statistical Society, Series B*, 68, 2, 239–257.
8
9 5 Statistics Canada (2007) 2005 Survey of Financial Security. Public Use Microdata File User
10 6 Guide. Ottawa: Statistics Canada. (Available from
11 7 <http://www.statcan.gc.ca/pub/13f0026m/13f0026m2007001-eng.htm>.)
12
13 8 Wolter K.M. (1985). *Introduction to variance estimator*, New York, Springer-Verlag.
14
15 9 You, Y. & Rao, J.N.K. (2002). Small area estimation using unmatched sampling and linking
16 10 models. *Canadian Journal of Statistics*, 30, 3–15.
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 **Appendix 1**

4
5
6 **Proof of [6]**

7 To start with we note that

8
9
10
$$\hat{\theta}_d^{B1} = \exp\left\{\mathbf{x}'_d \boldsymbol{\beta} + \gamma_d (\hat{Z}_d - \mathbf{x}'_d \boldsymbol{\beta}) + \frac{1}{2} \gamma_d \hat{\delta}_d\right\} = \exp\left\{\gamma_d \hat{Z}_d + (1 - \gamma_d) \mathbf{x}'_d \boldsymbol{\beta} + \frac{1}{2} \gamma_d \hat{\delta}_d\right\}$$
 so

11
12
13
$$E\left[\exp\left\{\gamma_d \hat{Z}_d + (1 - \gamma_d) \mathbf{x}'_d \boldsymbol{\beta} + \frac{1}{2} \gamma_d \hat{\delta}_d\right\}\right] = \exp\left\{(1 - \gamma_d) \mathbf{x}'_d \boldsymbol{\beta} + \frac{1}{2} \gamma_d \hat{\delta}_d\right\} E\left[\exp\{\gamma_d \hat{Z}_d\}\right].$$

14
15
16 We note that $E(\hat{Z}_d) = E_M\{E(\hat{Z}_d | \eta_d)\} = E_M\{\eta_d\} = \mathbf{x}'_d \boldsymbol{\beta}$; analogously

17
18
19
$$V(\hat{Z}_d) = V_M\{E(\hat{Z}_d | \eta_d)\} + E_M\{V(\hat{Z}_d | \eta_d)\} = \boldsymbol{\psi}_d + \hat{\delta}_d$$
 . As a consequence $E(\gamma_d \hat{Z}_d) = \gamma_d \mathbf{x}'_d \boldsymbol{\beta}$,

20
21
22
$$V(\gamma_d \hat{Z}_d) = \gamma_d^2 (\boldsymbol{\psi}_d + \hat{\delta}_d) = \boldsymbol{\psi}_d \gamma_d$$
 and $E\left[\exp\{\gamma_d \hat{Z}_d\}\right] = \exp\left\{\gamma_d \mathbf{x}'_d \boldsymbol{\beta} + \frac{1}{2} \boldsymbol{\psi}_d \gamma_d\right\}$. This leads to

23
24
25
$$E(\hat{\theta}_d^{B1}) = \exp\left\{\mathbf{x}'_d \boldsymbol{\beta} + \frac{1}{2} \gamma_d (\boldsymbol{\psi}_d + \hat{\delta}_d)\right\} = \exp\left\{\mathbf{x}'_d \boldsymbol{\beta} + \frac{1}{2} \boldsymbol{\psi}_d\right\}$$
 that coincides with $E_M(\theta_d)$.

26
27
28
29
30
31

32 **Proof of [7]**

33 Let's introduce some matrix notation. Let $\mathbf{z} = \text{vec}(\hat{Z}_d)$ is the vector containing, $\boldsymbol{\Psi} = \text{diag}(\boldsymbol{\psi}_d)$,

34 $\boldsymbol{\Lambda} = \text{diag}(\hat{\delta}_d)$ the diagonal matrixes containing the variance components; let \mathbf{X} be the matrix with

35
36
37
38 rows \mathbf{x}'_d , $d = 1, \dots, m$.

39
40
41 Standard Bayesian analysis of normal linear mixed model lead to $\boldsymbol{\beta} | \mathbf{z}, \boldsymbol{\Psi} \sim N(\hat{\boldsymbol{\beta}}_{gls}, \mathbf{V}(\boldsymbol{\Psi}))$ where

42
43
$$\hat{\boldsymbol{\beta}}_{gls} = \left[\mathbf{X}'(\boldsymbol{\Psi} + \boldsymbol{\Lambda})^{-1} \mathbf{X}\right]^{-1} \mathbf{X}'(\boldsymbol{\Psi} + \boldsymbol{\Lambda})^{-1} \mathbf{z}$$
 , $\mathbf{V}(\boldsymbol{\Psi}) = \left[\mathbf{X}'(\boldsymbol{\Psi} + \boldsymbol{\Lambda})^{-1} \mathbf{X}\right]^{-1}$. We can calculate

44
45
$$E(\boldsymbol{\eta} | \mathbf{z}, \boldsymbol{\Psi}) = E_{\boldsymbol{\beta} | \mathbf{z}, \boldsymbol{\Psi}} E(\boldsymbol{\eta} | \mathbf{z}, \boldsymbol{\beta}, \boldsymbol{\Psi})$$
 where $\boldsymbol{\eta} = \text{vec}(\eta_d)$. $E_{\boldsymbol{\beta} | \mathbf{z}, \boldsymbol{\Psi}} E(\boldsymbol{\eta} | \mathbf{z}, \boldsymbol{\beta}, \boldsymbol{\Psi}) = E_{\boldsymbol{\beta} | \mathbf{z}, \boldsymbol{\Psi}} (\mathbf{X} \boldsymbol{\beta}) = \mathbf{X} \hat{\boldsymbol{\beta}}_{gls}$.

46
47 Analogously $V(\boldsymbol{\eta} | \mathbf{z}, \boldsymbol{\Psi}) = V_{\boldsymbol{\beta} | \mathbf{z}, \boldsymbol{\Psi}} E(\boldsymbol{\eta} | \mathbf{z}, \boldsymbol{\beta}, \boldsymbol{\Psi}) + E_{\boldsymbol{\beta} | \mathbf{z}, \boldsymbol{\Psi}} V(\boldsymbol{\eta} | \mathbf{z}, \boldsymbol{\beta}, \boldsymbol{\Psi})$. If we denote the vector of small

48
49 area predictors (on the log scale) conditional on $\boldsymbol{\beta}$ as $\boldsymbol{\eta}^{B1} = E(\boldsymbol{\eta} | \mathbf{z}, \boldsymbol{\beta}, \boldsymbol{\Psi}) = \boldsymbol{\Gamma} \mathbf{z} + (\mathbf{I} - \boldsymbol{\Gamma}) \mathbf{X} \boldsymbol{\beta}$ with

50
51
$$\boldsymbol{\Gamma} = \boldsymbol{\Psi}(\boldsymbol{\Psi} + \boldsymbol{\Lambda})^{-1}$$
 we have that $V(\boldsymbol{\eta} | \mathbf{z}, \boldsymbol{\Psi}) = V_{\boldsymbol{\beta} | \mathbf{z}, \boldsymbol{\Psi}} \boldsymbol{\Gamma} \mathbf{z} + (\mathbf{I} - \boldsymbol{\Gamma}) \mathbf{X} \boldsymbol{\beta} + E_{\boldsymbol{\beta} | \mathbf{z}, \boldsymbol{\Psi}} \mathbf{G}_1$ with $\mathbf{G}_1 = \boldsymbol{\Psi}(\boldsymbol{\Psi} + \boldsymbol{\Lambda})^{-1} \boldsymbol{\Lambda}$

52
53
54 . Taking expectation with respect to $p(\boldsymbol{\beta} | \mathbf{z}, \boldsymbol{\Psi})$ we get

1 $V(\boldsymbol{\eta} | \mathbf{z}, \boldsymbol{\Psi}) = (\mathbf{I} - \boldsymbol{\Gamma}) \mathbf{X}' \boldsymbol{\beta} \mathbf{V}(\boldsymbol{\Psi}) \mathbf{X} (\mathbf{I} - \boldsymbol{\Gamma}) + \mathbf{G}_1 = \mathbf{G}_2 + \mathbf{G}_1$. We note that $p(\boldsymbol{\eta} | \mathbf{z}, \boldsymbol{\Psi})$ is a multivariate

2 normal distribution. If we consider an individual η_d have that $E(\eta_{ijr} | \mathbf{z}, \boldsymbol{\Psi}) = \mathbf{x}_d' \hat{\boldsymbol{\beta}}_{gls}$,

3
$$V(\eta_d | \mathbf{z}, \boldsymbol{\Psi}) = \gamma_d \hat{\delta}_d + (1 - \gamma_d)^2 \mathbf{x}_d' \left(\sum_{d=1}^m \frac{1}{\psi_d + \hat{\delta}_d} \mathbf{x}_d \mathbf{x}_d' \right)^{-1} \mathbf{x}_d = g_{1,d} + g_{2,d}.$$

4 As $\theta_d = \exp(\eta_d)$ formula [7] follows.

5

6

7 **Proof of [11] and subsequent statements**

8 From [9] we have that $p(\boldsymbol{\Psi} | \lambda) = \prod_{d=1}^m \frac{\lambda^a}{\Gamma(a)} \psi_d^{a-1} e^{-\lambda \sum_{d=1}^m \psi_d} \propto \lambda^{ma} e^{-\lambda \sum_{d=1}^m \psi_d} \prod_{d=1}^m \psi_d^{a-1}$

9 Conditioning on a is omitted as it is treated as a known constant; m is a shortcut notation for the
10 overall number of areas.

11 We can obtain the marginal prior $p(\boldsymbol{\Psi})$ using the integral $p(\boldsymbol{\Psi}) = \int_{\mathbb{R}^+} p(\boldsymbol{\Psi} | \lambda) p(\lambda) d\lambda$. As

12 $\lambda \sim \text{Gamma}(b_0, c_0)$ we straightforwardly get

13
$$p(\boldsymbol{\Psi}) \propto \frac{\Gamma(ma + b_0)}{\left(\sum_{d=1}^m \psi_d + c_0 \right)^{ma + b_0}} \prod_{d=1}^m \psi_d^{a-1}$$

14 Applying the transformation $\tau_d = \sqrt{\psi_d}$ on each component of $\boldsymbol{\Psi}$ we obtain

15
$$p(\boldsymbol{\tau}) \propto \left(\sum_{d=1}^m \tau_{ijr}^2 + c_0 \right)^{-(ma + b_0)} \prod_{d=1}^m \tau_d^{2a-1}$$

16 For the special case $a = \frac{1}{2}$ the density of $p(\boldsymbol{\tau})$ simplifies to $p(\boldsymbol{\tau}) \propto \left(\sum_{d=1}^m \tau_d^2 + c_0 \right)^{-(ma + b_0)}$ or

17 equivalently to

18
$$p(\boldsymbol{\tau}) \propto \left(1 + \frac{1}{c_0} \sum_{d=1}^m \tau_d^2 \right)^{-\left(\frac{m}{2} + b_0 \right)},$$

19 $\tau_d > 0, \forall d$. This expression can be recognized as the kernel of the density of a multivariate half-t
20 distribution with mean vector 0 and diagonal scale matrix. A multivariate half-t is a multivariate t
21 for which we apply the absolute value transformation on each component. We can then write
22 formula [13].

23

1 To prove that univariate priors $p(\tau_d)$ let's start from $p(\tau_d) = \int \dots \int p(\boldsymbol{\tau}) d\boldsymbol{\tau}_{-d}$. We can represent as

2 $p(\boldsymbol{\tau})$ as the result of applying the absolute value transformation on a multivariate t distribution, i.e

3 $\boldsymbol{\tau} = |\boldsymbol{\tau}^*|$ with $p(\boldsymbol{\tau}^*) = \int p(\boldsymbol{\tau}^* | \xi) p(\xi) d\xi$ where $p(\boldsymbol{\tau}^* | \xi) \sim MVN\left(\mathbf{0}, \frac{2b_0}{c_0} \xi \mathbf{I}_m\right)$ and

4 $p(\xi) \sim IGamma(b_0, b_0)$. We use the fact that a random vector distributed according to a
5 multivariate t distribution can be represented as an inverse Gamma mixture of a multivariate
6 normal.

7 As the variance covariance matrix of $\boldsymbol{\tau}^*$ is diagonal $p(\boldsymbol{\tau}^*) = \int \prod_{d=1}^m p(\tau_d^* | \xi) p(\xi) d\xi$.

8 Horrace (2005) studies truncated multivariate normal distributions and shows that univariate
9 marginal distributions from a multivariate half-Normal (obtained applying the absolute value
10 transformation on each component) are univariate half-Normals if and only the variance-covariance
11 matrix of the parent multivariate normal is diagonal. As a consequence

$$12 \quad p(\tau_d) = \int \dots \int \left[\int \prod_{d=1}^m p(\tau_d | \xi) p(\xi) d\xi \right] d\boldsymbol{\tau}_{-d}$$

13 where each $p(\tau_d | \xi)$ is distributed as an half-Normal distribution.

14 If we change the order of integration and use conditional independence of $p(\tau_d | \xi)$ we obtain that

15 $p(\tau_d)$ are marginally half-t distributed.

16 To prove that τ_d are linearly independent of each other we write

$$17 \quad V(\boldsymbol{\tau}) = E_{\xi} \{V(\boldsymbol{\tau} | \xi)\} + V_{\xi} \{E(\boldsymbol{\tau} | \xi)\}$$

18 and note that $E(\boldsymbol{\tau} | \xi) = \mathbf{0}$ while $E_{\xi} \{V(\boldsymbol{\tau} | \xi)\} = E(\xi) \frac{2b_0}{c_0} \mathbf{I}_m$, which is of course diagonal.

Bayesian small area estimation for skewed business survey variables

Enrico Fabrizi

Università Cattolica del S. Cuore, Piacenza, Italy – enrico.fabrizi@unicatt.it

Maria Rosaria Ferrante, Carlo Trivisano

Università di Bologna, Bologna, Italy – maria.ferrante@unibo.it, carlo.trivisano@unibo.it

Abstract

In business surveys, estimates of means and totals for sub-national regions, industries and business classes can be too imprecise due to the small sample sizes available for sub-populations. We propose a small area technique for the estimation of totals for skewed target variables, which are typical of business data. We adopt a Bayesian approach to inference. We specify a prior distribution for the random effects based on the idea of local shrinkage, which is suitable when auxiliary variables with strong predictive power are available, another feature often displayed by business survey data. This flexible modelling of random effects lead to predictions in agreement with those based on global shrinkage for most of the areas, but allows to obtain less shrunken and thereby less biased estimates for areas characterized by large model residuals. We discuss an application based on data from the Italian survey on Small and Medium Enterprises. By means of a simulation exercise, we explore the frequentist properties of the proposed estimators. They are good, and with difference to methods based on global shrinkage remain so also for areas characterized by large model residuals.

Keywords: robust estimation, Log-Normal distribution, local shrinkage priors, regional studies, Variance Gamma distribution.

1. Introduction

Regional economic decisions and policies rely on accurate business information regarding sub-national regions and business categories. The relevance of regional estimates of business

1
2
3 1 aggregates and the interest in regional disparities in terms firm competitiveness and productivity is
4 2 demonstrated by the growing number of scientific articles in this field (see Breinlich et al., 2014 for
5 3 a review).

6 4 Regional statistics are produced by the National Statistical Institutes, and governments use
7 5 them to coherently allocate funds (for examples of this, see OECD, 2013, Eurostat, 2011, Eurostat,
8 6 2015a). For instance, the (gross) value added, that is, the total value of new goods produced and
9 7 services provided in a given time period, is routinely estimated at the national and sub-national
10 8 levels. For the EU, Eurostat releases regional estimates of the value added at levels as detailed as
11 9 the EU NUTS 3 regions (following the Nomenclature of Territorial Units for Statistics, Eurostat,
12 10 2015b) and industries (NACE Rev. 2, 1 digit, following the Statistical classification of economic
13 11 activities in the European Community). Sub-national estimates of value added would be even more
14 12 informative if they were disaggregated both in terms of industry and firm size for the purpose of
15 13 measuring the relative contribution of an industry and of certain firm-size classes to the regional
16 14 economy. Unfortunately, sample sizes of official business surveys are too small for the standard
17 15 design-based estimators (known as “direct estimators”) to be sufficiently precise in small domains.

18 16 This limitation can be overcome by model-based small area estimation methods. The small
19 17 area estimation literature has until very recently focused largely on the analysis of social surveys,
20 18 with estimation goals such as the poverty mapping (see Pfeiffermann, 2014 and Pratesi, 2016 for a
21 19 review) and few applications for business statistics. In the last few years, awareness of this field of
22 20 application has grown (Burgard et al. 2014; Ferrante and Trivisano, 2010; Militino et al., 2015), as
23 21 well as the availability of reliable administrative archives for firms that can be used to obtain
24 22 auxiliary information.

25 23 Small area models may be broadly classified into *area level* and *unit level*. In area level
26 24 models, survey weighted (direct) estimates obtained for each domain are related with auxiliary
27 25 information at the same level of population disaggregation. In unit level models, the target variables
28 26 and unit level information on auxiliary variables are related at this micro level. Area level models
29 27 straightforwardly incorporate information on survey weights, leading to design-consistent
30 28 estimators whenever direct estimators are design-consistent (Rao, 2003, p. 117). Design consistency
31 29 is a general purpose form of protection against model failures, as it guarantees that, at least for large
32 30 domains, estimates make sense even if the assumed model completely fails. Area level modelling is
33 31 less demanding in terms of data disclosure and overcomes potential problems of record linkage
34 32 between the survey sample and the administrative archive. For these reasons, area level models will
35 33 be considered in this paper.

1
2
3 1 Many business survey variables are positive and positively skewed (Rivière, 2002), so
4 2 normality is not a tenable assumption in most of the cases. Log transformation can then be
5 3 introduced in order to apply normal linear mixed models on the log scale. Predictions on the
6 4 original data scale require back-transformation that is a potential source of bias. Positive skewness
7 5 of survey variables may cause estimators of means and totals to have non-normal (positively
8 6 skewed) sampling distributions, when calculated on small samples (see Fay and Herriot, 1979;
9 7 Karlberg, 2000). Literature on area level modelling on the log scale include Fay and Herriot (1979)
10 8 and Slud and Maiti (2006) that both consider an empirical Bayes approach to inference. In this
11 9 paper we propose a full Bayes approach, accounting for all sources of uncertainty, effectively
12 10 dealing with back-transformation bias and implementable with widely available MCMC software.

13 11 When predicting means or totals for business survey variables, strong covariates from
14 12 administrative archives are often available. For instance, in our application, aimed at predicting
15 13 gross value added at the domain level, we can exploit the knowledge of turnover for each firm in
16 14 the population. Area level totals of turnover are strongly correlated with those of value added.
17 15 Nonetheless, a minority of the areas will typically deviate from the relationship that characterize
18 16 most of the others. If we think of modelling in terms of mixed models, we have that random effects
19 17 would be needed for a subset of the areas (Datta, 2011) or alternatively that there are subsets of
20 18 random effects characterized by different variances. The specification of spike-and-slab priors can
21 19 be useful in this case (Datta and Mandal, 2015).

22 20 We contribute to the small area literature by proposing an approach based on local shrinkage
23 21 priors for the random effects (Frühwirth-Schnatter and Wagner, 2010) where spike-and-slab priors
24 22 are replaced by continuous gamma scale mixture of normal distributions (Griffin and Brown, 2010)
25 23 that lead to marginal variance-gamma distributions for the random effects. This flexible modelling
26 24 of random effects lead to predictions close to those we can obtain using standard priors for non-
27 25 outlying areas, and to less biased predictors for the areas that can be labelled as outliers.

28 26 The paper is organized as follows. Model specification is described in section 2. Specifically
29 27 in section 2.1 closed formulas for posterior means conditionally on variance components are
30 28 illustrated as posterior means are proposed as point predictors. In section 3, we apply our
31 29 methodology on real survey data. We use data on the Small and Medium Enterprises sample survey
32 30 (1-99 employees) conducted by the Italian National Statistical Institute (ISTAT), which provided us
33 31 with this information within the framework of the BLUE-ETS project; this project has been
34 32 financially supported by the EU Commission within the 7th Framework Programme. For these data
35 33 we motivate the recourse to log normal likelihood for the direct estimators. In section 4, we
36 34 introduce a simulation exercise to explore the frequentist properties of the proposed predictor in

1 comparison with some alternatives, including the estimator of Slud and Maiti (2006). Section 5
2 presents the study's conclusions.

3 2. Small Area Estimation Model

4 Let Y be the target variable, which we assume positive with a positively skewed distribution.
5 Assume that Y is defined on a population U of N units, partitioned into a set of m non-overlapping
6 domains of size N_d ($d = 1, \dots, m; N = \sum_{d=1}^m N_d$). A random sample of overall size n is taken using a
7 possibly complex design: samples of sizes n_d are drawn from each domain. The small area nature
8 of the problem lies in n_d being too small to allow for reliable inference for most of the domains.
9 We assume that individual weights w_{dj} , $j = 1, \dots, n_d$ are attached to responses y_{dj} to account for
10 unequal selection probabilities and possibly other selection adjustments.

11 The normal distribution is not suitable to describe either the distribution of Y in the
12 population nor the sampling distribution of the domain totals' direct estimators $\hat{Y}_d = \sum_{j=1}^{n_d} w_{dj} y_{dj}$.
13 Although these are linear combinations of individual observations and can be assumed to be
14 approximately normally distributed in large samples, in samples of small size, the sum of a few
15 positively skewed variables remains positively skewed. We assume that the total direct estimators
16 are log-normally distributed:

$$17 \hat{Y}_d | \theta_d, V_d \sim LN([\theta_d], [V_d]) \quad [1]$$

18 where $[\cdot]$ is used to denote a parametrization in terms of mean and variance of the distribution.
19 Exact or approximate design-unbiasedness of totals' estimators is typical in survey sampling. The
20 distributional assumption in [1] can be motivated directly assuming the log-normality of Y . Log-
21 normal approximations of sums of independent log-normals are justified by several authors (e.g.,
22 Fenton, 1960; Cobb et al., 2012). Moreover, Mazmanyany et al. (2009) proposed a log-normal
23 central limit theorem for the approximation of the sum of positively skewed random variables,
24 although not necessarily log-normal. Eventually, the assumption of normality on the log scale when
25 dealing with mean or total estimators of skewed variables is common in the small area literature (as
26 in Fay and Herriot, 1979).

27 On the log-scale, a specification consistent with the *sampling model* [1] is given by:

$$28 \log(\hat{Y}_d) | \eta_d, \delta_d \sim N(\eta_d - \delta_d/2, \delta_d) \quad [2]$$

1 where $\eta_d = \ln(\theta_d)$ and $\delta_d = \text{Var}\{\log(\hat{Y}_d)\}$. $E\{\log(\hat{Y}_d)\} = \eta_d - \delta_d/2$ is in line with assuming the
 2 availability of an unbiased estimator on the original scale of the data: if $E(\hat{Y}_d) = \theta_d$, then
 3 $E\{\log(\hat{Y}_d)\} < \log(\theta_d)$. Note that $V_d = \exp\{2\theta_d + \delta_d\}[\exp\{\delta_d\} - 1]$ will depend on both parameters
 4 of the lognormal distribution.

5 In the small area literature, variances associated with direct estimators are usually treated as
 6 known constants. In practice, estimates obtained with methods such as linearization or bootstrap are
 7 smoothed using a model involving unknown parameters. In line with the literature on area-level
 8 models, we will assume that variances on the log-scale are known and denote them as $\hat{\delta}_d$.

9 We assume a multiplicative *linking model* for θ_d that links the outcome parameter to the
 10 auxiliary information in order to improve the direct estimators:

$$11 \quad \theta_d = \exp(\mathbf{x}'_d \boldsymbol{\beta} + u_d) \quad [3]$$

12 The p -row vector \mathbf{x}'_d contains the covariates known for domain d from external sources, and u_d is a
 13 random intercept associated to θ_d . Let us assume that $u_d \sim N(0, \psi_d)$, which implies

$$14 \quad \theta_d \sim LN(\mathbf{x}'_d \boldsymbol{\beta}, \psi_d) \quad [4]$$

15 or, equivalently, $\eta_d \sim N(\mathbf{x}'_d \boldsymbol{\beta}, \psi_d)$. We denote the model defined by *sampling model* [1] and *linking*
 16 [4] as the LN-LN model.

19 2.1 Analysis conditional on the variance components

20 To analyze the model defined by [1] and [4], note first that, assuming δ_d as known
 21 ($\delta_d = \hat{\delta}_d$) we can re-write [2] as $\hat{Z}_d \sim N(\eta_d, \hat{\delta}_d)$, where $\hat{Z}_d = \log \hat{Y}_d + \frac{1}{2} \hat{\delta}_d$. We can use standard
 22 results from the analysis of linear mixed models (see Rao 2003, chapter 5) to prove that,
 23 conditionally on the regression coefficients $\boldsymbol{\beta}$ and the variances ψ_d :

$$24 \quad \eta_d | \boldsymbol{\beta}, \psi_d, \text{data} \sim N(\hat{\eta}_d^{B1}, g_{1,d})$$

25 where $\hat{\eta}_d^{B1} = \mathbf{x}'_d \boldsymbol{\beta} + \gamma_d (\hat{Z}_d - \mathbf{x}'_d \boldsymbol{\beta})$, $g_{1,d} = \gamma_d \hat{\delta}_d$ and $\gamma_d = \frac{\psi_d}{\psi_d + \hat{\delta}_d}$. Note that, as a function of the

26 shifted direct estimates \hat{Z}_d , $\hat{\eta}_d^{B1}$ is a convex linear combination of a direct (\hat{Z}_d) and a synthetic
 27 component ($\mathbf{x}'_d \boldsymbol{\beta}$), known as the linear composite estimator in the small area literature. If we

1 assume quadratic loss and define $\hat{\theta}_d^{B1} = E(\theta_d | \boldsymbol{\beta}, \boldsymbol{\psi}_d, \text{data})$ as the point predictor for θ_d , we have
 2 that

$$\begin{aligned} \hat{\theta}_d^{B1} &= \exp\left\{\mathbf{x}'_d \boldsymbol{\beta} + \gamma_d (\hat{Z}_d - \mathbf{x}'_d \boldsymbol{\beta}) + \frac{1}{2} \gamma_d \hat{\delta}_d\right\} \\ &= \exp\left\{\mathbf{x}'_d \boldsymbol{\beta} + \gamma_d (\log \hat{Y}_d - \mathbf{x}'_d \boldsymbol{\beta}) + \gamma_d \hat{\delta}_d\right\} \end{aligned} \quad [5]$$

3 This predictor is the product between $\exp(\hat{\eta}_d^{B1})$ and a factor that corrects for the main bias term in
 4 the back-transformation; it is in line with formula (4) of Slud and Maiti (2006).

5 It can also be shown that

$$E(\hat{\theta}_d^{B1}) = \exp\left\{\mathbf{x}'_d \boldsymbol{\beta} + \frac{1}{2} \boldsymbol{\psi}_d\right\} = E_M(\theta_d) \quad [6]$$

7 where $E(\hat{\theta}_d^{B1})$ is the expectation taken with respect to both linking and sampling models, while
 8 with $E_M(\cdot)$, we denote the expectation with respect to linking model [4]. This latter result means
 9 that $\hat{\theta}_d^{B1}$ is an unbiased predictor of θ_d in the same sense that BLUP predictors are unbiased: the
 10 unconditional frequentist expectation of the estimator and the expectation of the estimand under the
 11 linking model are the same. A proof of [6] can be found in Appendix 1.

12 If we remove the conditioning on $\boldsymbol{\beta}$ and assume a non-informative flat prior on $\boldsymbol{\beta}$, i.e.,

13 $p(\boldsymbol{\beta}) \propto \mathbf{1}$, then we have that

$$\eta_d | \boldsymbol{\psi}_d, \text{data} \sim N(\hat{\eta}_d^{B2}, g_{1,d} + g_{2,d})$$

14 where $\hat{\eta}_d^{B2} = \mathbf{x}'_d \hat{\boldsymbol{\beta}}_{gls} + \gamma_d (\hat{Z}_d - \mathbf{x}'_d \hat{\boldsymbol{\beta}}_{gls})$, $\hat{\boldsymbol{\beta}}_{gls} = \left(\sum_d \frac{1}{\boldsymbol{\psi}_d + \hat{\delta}_d} \mathbf{x}_d \mathbf{x}'_d\right)^{-1} \sum_d \frac{1}{\boldsymbol{\psi}_d + \hat{\delta}_d} \mathbf{x}'_d \log \hat{Y}_d$,

15 $g_{2,d} = (1 - \gamma_d)^2 \mathbf{x}'_d \left(\sum_d \frac{1}{\boldsymbol{\psi}_d + \hat{\delta}_d} \mathbf{x}_d \mathbf{x}'_d\right)^{-1} \mathbf{x}_d$. As a consequence, the point predictor under quadratic

16 loss will be given by

$$\hat{\theta}_d^{B2} = \exp\left\{\mathbf{x}'_d \hat{\boldsymbol{\beta}}_{gls} + \gamma_d (\hat{Z}_d - \mathbf{x}'_d \hat{\boldsymbol{\beta}}_{gls}) + \frac{1}{2} (g_{1,d} + g_{2,d})\right\} \quad [7]$$

17 (see Appendix 1 for a proof). Unlike the empirical Bayes approach advocated by Slud and Maiti
 18 (2006), who plug estimates of unknown parameters into [7], a full Bayes approach accounts for the
 19 effect that the extra-variation implied by the estimation of $\boldsymbol{\beta}$ has on the point predictor; in fact, the
 20 expectation of a log-normal variable depends on both the expectation and variance on the log scale.

1 To fully account for all sources of uncertainty, we should remove the conditioning on the
 2 variance components ψ_d ; unfortunately, for sensible choices of the prior, this leads to posterior
 3 distributions for θ_d that cannot be written in closed form and should therefore be explored by
 4 means of computational algorithms such as the Markov Chain Monte Carlo considered in this
 5 paper.

6 2.2 The distribution for the random effects and specification of hyperpriors

7 The main difference between [4] and the linking model adopted by most of the small area
 8 literature on Fay-Herriot type models (Jiang and Lahiri, 2005; Pfeffermann, 2014) is that the
 9 variances associated with random intercepts are in [4] domain-specific, implying local shrinkage
 10 instead of the ordinary global shrinkage that we would have had assuming $\psi_d = \psi, \forall d$. In a
 11 different context, the specification of a distribution for random intercept based on local shrinkage is
 12 discussed in Frühwirth-Schnatter and Wagner (2010).

13 Datta et al. (2011) note that in the presence of good covariates, the variability of the small
 14 area parameters may be accounted for by a synthetic estimator, and the inclusion of a random effect
 15 term may be unnecessary. When random effects are needed for a subset of the areas, the
 16 specification of spike-and-slab priors can be useful (see Datta and Mandal, 2015). Spike-and-slab
 17 priors amount to assuming that random intercepts are sampled from a mixture of two normal
 18 distributions.

19 When analyzing business data, it is quite likely that auxiliary variables with strong
 20 predictive power are available. When this is the case, the bulk of the direct estimates will be well
 21 fitted by the synthetic model (without random intercepts), so the associated ψ_d are likely to be
 22 small, with a minority of areas that require larger area-specific intercepts (and thereby larger ψ_d).

23 Our specification for the distribution of $u_d, d = 1, \dots, m$ is based on infinite mixtures of
 24 normal distributions. Following Griffin and Brown (2010), our specification uses *Gamma* mixtures
 25 of normal distributions. Specifically, we assume:

$$26 u_d | \psi_d \stackrel{ind}{\sim} N(0, \psi_d) \quad [8]$$

$$27 \psi_d | a, \lambda \stackrel{ind}{\sim} \text{Gamma}(a, \lambda) \quad [9]$$

$$28 \lambda | b_0, c_0 \sim \text{Gamma}(b_0, c_0) \quad [10]$$

29 This leads to spiked priors for the random effects that at the same time have tails that are
 30 heavier than those of the normal distribution. Griffin and Brown (2010) observe that for small

values of the shape parameter a , the prior specification [8]-[10] leads to a marginal prior distribution for u_d that mimics the behavior of spike-and-slab priors based on finite mixtures. This infinite mixture specification is computationally easier to deal with.

Other choices for the mixing distribution such as the popular *Inverse Gamma* would lead to platikurtic distributions with heavy tails, such as those in the t family; this contrasts with the idea of severe shrinkage for most of the areas, which is consistent with a large probability mass close to 0.

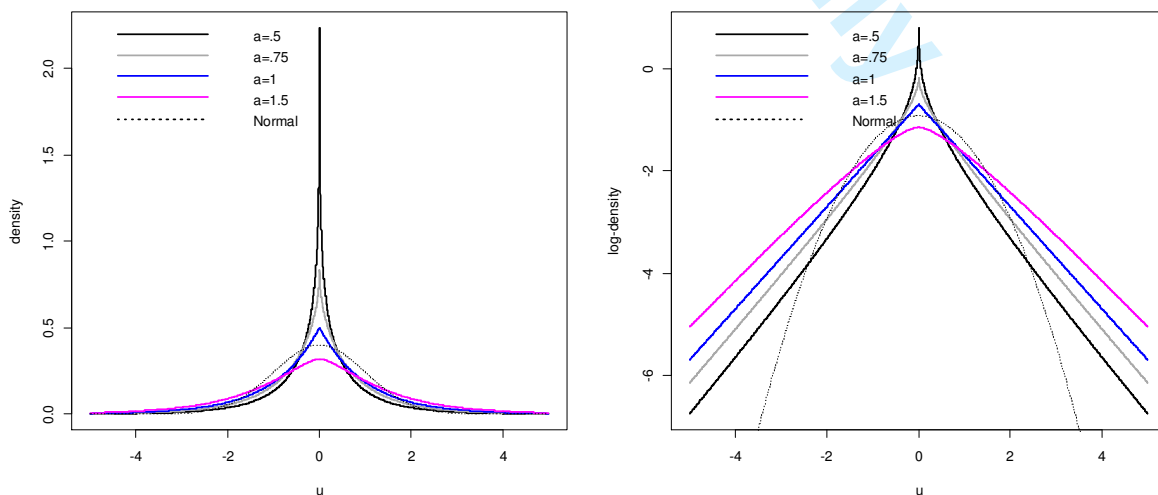
Specifically, prior specification [8]-[10] implies that $u_d|a,\lambda$ follows a Variance Gamma distribution, i.e.,

$$u_d \sim VG(a, \sqrt{2\lambda}, 0, 0)$$

(see Bibby and Sørensen, 2003 for more details on this distribution). This marginal prior distribution is symmetric and has expectation $E(u_d) = 0$ and variance $V(u_d) = a/\lambda$. It belongs to the family of generalized hyperbolic distribution (Barndorff-Nielsen, 1977). The conjugate hierarchy in [8]-[10] also facilitates MCMC sampling.

In line with Frühwirth-Schnatter and Wagner (2010), we set the shape parameter a to a fixed (small) value, while we specify a prior on the global parameter λ . As far as a is concerned, we focus on two choices, $a = 1$, $a = 0.5$.

The choice $a = 1$ implies that u_d is a priori distributed as a double-exponential or Laplace, which, combined with the normal prior conditional on ψ_d , recalls the Bayesian lasso of Park and Casella (2008).



21

1
2
3 1 Figure 1: *density (left) and log-density (right) functions of the Variance Gamma distribution.*

4 2 $VG(a, \sqrt{2\lambda}, 0, 0)$ for $\lambda = 1$ and various values of a .

5
6
7 3
8
9 4 The case $a = 0.5$ represents a more peaked prior distribution and encourages more shrinkage
10 5 towards 0 of small random intercepts (see Figure 1). Moreover, it leads to a *half-t* marginal prior on
11 6 $\sqrt{\psi_d}$. The *half-t* prior for standard deviations is discussed in Gelman (2006) and recommended
12 7 whenever it makes sense to put a sizeable mass of prior probability close to 0. It can be shown that
13 8 once $\boldsymbol{\psi} = \{\psi_d\}$ and $\boldsymbol{\tau} = \{\sqrt{\psi_d}\}$, $d = 1, \dots, m$, are defined, under prior [8]-[10] and $a = 0.5$,

14
15
16 9
$$\boldsymbol{\tau} | b_0, c_0 \sim Mht\left(0, \frac{2b_0}{c_0} \mathbf{I}, 2b_0\right). \quad [11]$$

17
18
19
20
21
22 10 With $Mht(\cdot)$ we denote the multivariate distribution (with support \mathbb{R}^{+n}) that is obtained from a
23 11 multivariate t distribution applying the absolute value transformation on each component of the
24 12 random variable. We can also prove that each prior $\sqrt{\psi_d} | b_0, c_0$ is a (univariate) *half-t* and the priors
25 13 for two different variance components are uncorrelated, i.e. $Cov(\psi_d, \psi_{d^*} | b_0, c_0) = 0$ whenever
26 14 $d \neq d^*$. See appendix 1 for a proof of [11] and the other statements.

27
28
29
30
31
32 15 As for the prior specification of the remaining parameters, diffuse independent normal priors
33 16 can be specified for the components of $\boldsymbol{\beta}$. We can set $b_0 = 2$, which implies that $E(\lambda^{-1}) = c_0$. This
34 17 helps to interpret c_0 as a scaling constant for the random effects variance $V(u_d) = a / \lambda$. The choice
35 18 of c_0 depends on the scale for the random effects in the problem being considered. According to
36 19 [11], the parameter b_0 can be interpreted in terms of degrees of freedom of the marginal prior
37 20 $p(\boldsymbol{\tau})$, so the choice of $b_0 = 2$ is in line with selecting *half-t* priors with a very small number of
38 21 degrees of freedom (Gelman, 2006).
39
40
41
42
43
44
45
46
47
48

49 24 **3. Estimation for Italian Small and Medium Enterprise Survey Data: An application**

50
51 25 In this section, we illustrate the methodology we discussed using real survey data. We use
52 26 data on the Small and Medium Enterprises (SME) sample survey, wave 2008, conducted by the
53 27 Italian National Statistical Institute (ISTAT). Specifically, we target the estimation of the total value
54 28 added (VA) for small domains of the population of Italian small and medium manufacturing firms
55 29 (less than 100 employees). The domains we focus on are smaller than those for which ISTAT
56 30 provides reliable estimates. Specifically, our domains are defined by cross-classifying: the 20 Italian

1 NUTS 2 administrative regions, the economic industrial sector (NACE Rev. 2, 2 digit, 22
2 industries) and firm size (4 classes: less than 10 employees, 10 to 19 employees, 20 to 49
3 employees, 50 to 99 employees). As anticipated, for domains as small as those that we target,
4 standard design-based estimators are characterized by unacceptably large variances.

3.1 Direct estimators and sampling model

7 The SME survey uses stratified sampling design and strata are defined by cross-classifying
8 NACE Rev. 2 (4 digits) Italian administrative regions and company size in the four classes defined
9 above. The domains we are interested in are planned because they are unions of sampling strata.

10 Let \hat{Y}_{ijr} be the direct estimator of the parameter θ_{ijr} , where i indexes the economic activity
11 ($i = 1, \dots, 22$), j the size classes ($j = 1, \dots, 4$) and r the regions ($r = 1, \dots, 20$). Given this peculiar
12 structure of the data the domain index d used in section 2 is now replaced by the triplet ijr . The
13 potential number of 1760 domains falls to 1165, as some of the populations obtained by cross-
14 classification are empty and some very small. We excluded all the domains characterized by a
15 sampling rate over 0.75.

16 The actual sample sizes for the domains we consider ranges from 2 to 184, with a median of
17 8, a mean of 13.5 and .75 and .9 quantiles equal to 16 and 30 respectively.

18 Direct estimates can be obtained using the calibration estimator that ISTAT adopts for the
19 SME survey. Calibration estimators can be written as weighted sums. ISTAT's published weights
20 are obtained by multiplying base weights (the inverse of the inclusion probabilities) by factors
21 adjusting for non-response and calibrating to known totals. Let the estimated total be denoted as

$$\hat{Y}_{ijr} = \sum_{k \in d_{ijr}} w_{ijr,k} y_{ijr,k}, \text{ where } y_{ijr,k} \text{ is the value added of the } k\text{-th firm in sector } i, \text{ size class } j, \text{ region } r.$$

23 We assume that $E(\hat{Y}_{ijr}) = \theta_{ijr}$ with $Var(\hat{Y}_{ijr}) = V_{ijr}$. We estimate V_{ijr} using linearization-based
24 variance estimators and denote these estimates as \hat{V}_{ijr} .

25 In our sampling model we assume log-normality according to [1]. To justify this assumption
26 for our data we proceed in two steps: first, we check whether log-normality is a sensible assumption
27 for domain specific sample data, then we use a simple simulation exercise to assess whether log-
28 normality is to be preferred to normality as the sampling distribution of total estimators given the
29 sample sizes we have in our analysis.

30 For all the domains with $n_{ijr} \geq 3$, we tested normality and log-normality using the Shapiro-
31 Wilk test. Results are reported in table 1 below. In reporting the results we consider separately the

1 smallest 90% of the domains $n_{ijr} \leq 30$ and the largest 10%. In the smaller domains, for which the
 2 test is relatively less powerful, both normality and log-normality tend to be not rejected, but
 3 normality fails clearly more often. In larger domains, when the test has more power, normality is
 4 rejected in the large majority of cases, while log-normality is accepted in more than 70% of the
 5 cases.

6
 7 Table 1: *Checking normality and log-normality within domain-specific samples using Shapiro-Wilk*
 8 *test. Percentage of non-rejections at the 0.01 significance level*

n_{ijr}	normality	log-normality
≤ 30	0.733	0.959
> 30	0.087	0.713
overall	0.672	0.943

9
 10 From table 1 we conclude that log-normality is a sensible assumption for the distribution of
 11 VA within domains. We actually assume that direct estimators are log-normally distributed
 12 according to the arguments illustrated in previous section. To check this, we consider a set of log-
 13 normal populations: $Y_d \sim LN(\tilde{\mu}_d, \tilde{\sigma}_d)$, $d = 1, \dots, L$ where $L = 77$ is the number of domains with
 14 $n_{ijr} > 30$ for which log-normality was not rejected and $\tilde{\mu}_d, \tilde{\sigma}_d$ are the parameters according to
 15 maximum likelihood for these domains. For each of these populations we generated simple random
 16 samples of sizes $R = 10,000$ for each of the following sample sizes: $n_d = 5, 10, 15, 20$. Note that 20
 17 represent the 0.8 quantile of the sample size distribution in our application.

18 We evaluate how far is the empirical sampling distribution of the sample mean from the
 19 normal and the log-normal distributions in terms of Kolmogorov-Smirnov distance averaging over
 20 the L populations. In fact, formal hypothesis testing of distributional assumptions with a sample of
 21 replicates as large as 10,000 would lead to rejections even in presence of negligible departures from
 22 the null. Results, summarized in table 2, show how log-normality is to be preferred to normality for
 23 all sample sizes. We can also note that as the sample size grows larger, the difference between the
 24 two distances decreases.

1
2
3 Table 2: *Kolmogorov-Smirnov distances between the Monte Carlo ($R=10,000$) distribution of the*
4
5 *sample mean and two reference distributions, for various sample sizes.*
6

	sample size (n_i)			
Reference distribution	5	10	15	20
Log-normal	0.012	0.014	0.015	0.015
Normal	0.119	0.097	0.081	0.073

7
8
9
10
11
12
13
14
15
16
17 In order to obtain more stable direct variance estimates, we smooth them through the
18 Generalized Variance Function approach (Wolter, 1986). To begin with, we consider that under the
19 log-normality assumption introduced in [1], we have that

$$20 \text{Var}\left\{\log\left(\hat{Y}_{ijr}\right)\right\} = \log\left\{CV^2\left(\hat{Y}_{ijr}\right) + 1\right\}. \quad [12]$$

21
22
23
24
25
26 Thus, the smoothing can be conducted on $CV^2\left(\hat{Y}_{ijr}\right) = \frac{\hat{V}_{ijr}}{\hat{Y}_{ijr}^2}$. After careful explorative analysis, we
27
28
29 assume that $CV^2\left(\hat{Y}_{ijr}\right)$ varies with the size class (j) but not with economic activity (i) or with
30
31 regions (r). This leads to the following smoothing equation for the direct estimate of V_{ijr} , \hat{V}_{ijr} :

$$32 \hat{V}_{ijr} = \phi_j \frac{\hat{Y}_{ijr}^2}{n_{ijr}} \left(1 - \frac{n_{ijr}}{N_{ijr}}\right) + v_{ijr} \quad [13]$$

33
34
35
36
37 with $E(v_{ijr}) = 0$, $V(v_{ijr}) = \kappa$ and where a finite population correction factor is also considered to
38
39 account for varying and occasionally non-negligible sample rates. The parameter ϕ_j can be
40
41 interpreted as the smoothed squared coefficient of variation multiplied for the size of the domain
42
43 n_{ijr} . The domain sample size n_{ijr} in the denominator of [13] allows for the decrease in the
44
45 coefficient of variation when the sample size increases. Smoothed squared estimated coefficient of
46
47 variations are given by $CV_{smooth}^2\left(\hat{Y}_{ijr,k}\right) = \frac{\phi_j}{n_{ijr}} \left(1 - \frac{n_{ijr}}{N_{ijr}}\right)$; the first, second and third quartiles of
48
49 $CV_{smooth}^2\left(\hat{Y}_{ijr,k}\right)$ estimated on our data set are 31%, 45% and 65%, respectively. These results
50
51 confirm the need to adopt a small area model approach.
52
53

54
55 We can then adapt the sampling model [2] to our problem changing the index from d to ijr
56
57 and δ_{ijr} with $\hat{\delta}_{ijr} = \log\left\{CV_{smooth}^2\left(\hat{Y}_{ijr}\right) + 1\right\}$ defined according to [12].
58
59
60

3.2 Auxiliary information and linking model

As an auxiliary variable, the log total turnover in each domain is available. This auxiliary information refers to the Italian firms' population and it is provided by the Italian Statistical Register of Active Enterprises-ASIA archive. The predictive power of this covariate is quite strong: the squared correlation coefficient is equal to 0.87 when calculated on variables on their original scale, and it is equal to 0.79 for the log transformations. In the original scale the high correlation level is influenced by few observations with a larger scale with respect to most of the others.

We assume the multiplicative *linking model* [4] for θ_{ijr} to link the outcome parameter to the auxiliary information given by the log-total turnover for the domain in question. With reference to log-scale, we can write $\eta_{ijr} = \beta_0 + \text{lt}_{ijr}\beta_1 + u_{ijr}$. The prior for the vector of domain-specific random intercepts u_{ijr} is specified according to [8]-[10]. As for the prior specifications not already discussed, we set $\beta_0 \sim N(0, 10^5)$, $\beta_1 \sim N(0, 10^5)$, $b_0 = 2$, $c_0 = 1$. We chose these values as they provide a reasonable scale for the random effects variance in our problem.

We also consider the LN-LN model with an alternative choice for the prior distribution on

u_{ijr} :

$$u_{ijr} | \sigma^2 \sim N(0, \sigma^2), \sigma^2 \sim \text{InverseGamma}(c, d) \quad [14]$$

This prior specification, which implements global shrinkage, can be considered as a benchmark for evaluating the effects of prior specification approximating spike-and-slab introduced in the previous section, and it represents a routine choice in many applications. We set $c = 0.01$, $d = 0.01$.

3.3 Markov Chain Monte Carlo computational issues

Parameter estimates are obtained by summarizing the posterior distributions approximated by the output of Markov Chain Monte Carlo (MCMC) integration via the Gibbs sampling algorithm. By assuming a quadratic loss, the posterior means are adopted as estimates of the area specific parameters. Posterior variances are used as a measure of uncertainty. To carefully assess the convergence, we run three parallel chains of 25,000 runs each, the starting point being drawn from an over-dispersed distribution. The convergence of the Gibbs sampler was monitored by visual inspection of the chains' plots and autocorrelation diagrams and by means of the potential scale reduction known as the Gelman-Rubin statistic (Carlin and Louis, 2000, ch. 5). Both models

1 displayed fast convergence; we discarded the first 5,000 iterations from each chain. To obtain
 2 estimates, we used the OpenBugs software package, which can be downloaded for free on the
 3 internet and it is open source.

3.4 Comparing alternative models

4
 5
 6 In order to choose among competing models, we compute the Deviance Information
 7 Criterion (DIC) and the logarithm of the pseudo-marginal likelihood (LPLM, Ibrahim et al., 2001).
 8 A model is preferred if it displays a lower DIC value. Table 3 reports the DIC results for the whole
 9 set of small area models estimated. DIC values show that, in line with expectations, the log-
 10 normality assumption at the sampling level performs better in terms of DIC with respect to the
 11 model assuming normality. The ordering of alternative models using LPLM is consistent with that
 12 obtained using the DIC. The adoption of the Variance Gamma for the random intercepts u_{ijr} leads
 13 to a further reduction in DIC with respect to the more common specification [14].

14
 15
 16 Table 3: Comparison of alternative assumptions on the distributions of the random effects

Shrinkage	Prior on random intercepts u_{ijr}	a	DIC	LPLM	median CVR
global	[14]	-	15340	-7846	0.391
local	[10]-[12]	1	15230	-7808	0.421
local	[10]-[12]	0.5	15220	-7798	0.455

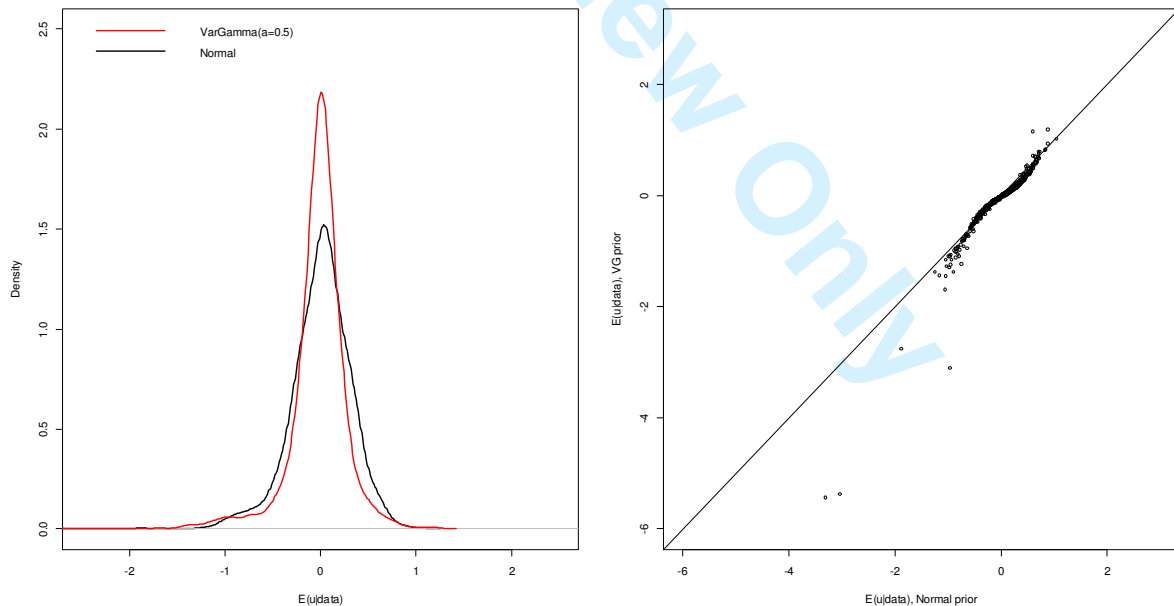
17
 18
 19 We also compare the median reduction of the coefficient of variation of estimators with
 20 respect to the direct ones, defined as $median(CVR_k^h)$. CVR_k^h is defined as $CVR_k^h = 1 - CV_k^h / CV_k^{DIR}$,
 21 where the CV_k^h is the coefficient of variation calculated on the posterior of θ_k (k being a generic
 22 index for the areas) under model h , while CV_k^{DIR} is the coefficient of variation of the direct
 23 estimators calculated from the randomization distribution.

24 The posterior predictive approach can be used to assess the fit of a model (Gelman et al.,
 25 1996). We consider a discrepancy measure suggested in the context of small area estimation by You
 26 and Rao (2002) and considered also in Fabrizi and Trivisano (2016):

$$27 \quad \text{dis}_{ijr} = P(\hat{Y}_{ijr} < Y_{ijr}^*)$$

1 where Y_{ijr}^* is generated from the posterior predictive distribution. The discrepancy measure is aimed
 2 at checking whether, for each area, the posterior predictive distribution is centered around the direct
 3 estimate. Values of dis_{ijr} far from 0 and 1 would provide evidence of systematic under or over
 4 estimation. For the log-normal model endowed with priors [10]-[12] and $a=0.5$ (i.e. the best
 5 model in terms of DIC and LPLM), we have that the average of the discrepancy measure over the
 6 set of areas is 0.499 with 0.25 and 0.75 quantiles equal 0.32 and 0.68 respectively, which means an
 7 adequate fit. Less than 10% of the areas shows dis_{ijr} out of the range (0.2, 0.8). Similar summaries
 8 can be obtained for the other models considered in table 3.

9 Results on $median(CVR_k^h)$, reported in table 3, highlight that the whole set of considered
 10 small area estimators considerably reduce the variability of direct estimators, which is consistent
 11 with the availability of a strongly predictive auxiliary variable. Nonetheless, even if exploiting the
 12 same auxiliary information, the models perform differently, as the prior specification has a non-
 13 negligible impact. Prior specifications mimicking the spike and slab behavior allow for a further
 14 gain in efficiency with respect to priors ordinarily used in this type of analysis.



16 Figure 2: Left panel: kernel densities of $E(u_{ijr} | data)$ over the set of small areas under Normal and
 17 VG priors. Right panel: $E(u_{ijr} | data)$ under Normal prior vs $E(u_{ijr} | data)$ under VG prior. For VG
 18 priors, $a=0.5$ is assumed.
 19
 20

1
2
3 1 To evaluate the improvements allowed by the model-based proposed predictor we can
4 2 compare the number of small areas with values of the coefficient of variation CV less than 16.6%,
5 3 between 16.6% and 33.3% and over 33.3% for the direct and the model-based predictor. These
6 4 thresholds for CV were suggested by Statistics Canada (2007) to provide quality level guidelines for
7 5 publishing small area estimates; those with a coefficient of variation less than 16.6% are considered
8 6 reliable for general use. Estimates with a coefficient of variation between 16.6% and 33.3% should
9 7 be accompanied by a warning to users. Estimates with coefficients of variation larger than 33.3%
10 8 are deemed to be unreliable. Less than 25% of the direct estimates have associated CV below the
11 9 33.3% threshold, while for the model based ones this number grows to 70%. Although the
12 10 uncertainty around the small area estimates remain sizeable and not all estimates would be
13 11 publishable, the application of the proposed method endows most sub-population with a publishable
14 12 estimate in spite of the small sample sizes.

15 13 Figure 2 displays the impact of alternative prior specifications on the ensemble of the
16 14 random intercepts' posterior means $E(u_{ijr} | data)$ under the Normal prior [14] and Variance Gamma
17 15 priors $u_{ijr} \sim VG(0.5, \sqrt{2\lambda}, 0, 0)$. Based on the left panel, it is clear that the Variance Gamma leads to
18 16 a more peaked distribution of estimated random intercepts, as predicted by theory. As tail behavior
19 17 is difficult to read from density estimates, in the right panel of Figure 2, we plotted the point
20 18 posterior expectation under the Normal prior vs those obtained under the Variance Gamma. The
21 19 peak around 0 is still apparent from the inflection of the points' cloud approximately at 0; heavier
22 20 tails under Variance Gamma can also be appreciated: under the Normal prior, $E(u_{ijr} | data)$ lies
23 21 within the interval $(-4, 2)$, while they do not under the Variance Gamma specification.

24 22 As the purpose of small area estimators is to complement the direct estimates obtained using
25 23 ordinary survey weighted methods, robustness with respect to modelling assumptions is a major
26 24 concern. As anticipated, the recourse to area-level model-based methods entails design consistency.
27 25 When the area-specific sample sizes are large (and occasionally they are) the small area estimate
28 26 will be close to that obtained under standard design-based methods. This offers, at least for the
29 27 larger domains, protection against model misspecifications; moreover, it automatically guarantees
30 28 that, in the case of large domains, model-based and design-based estimates are automatically in
31 29 agreement.

4. A simulation assessment of the frequentist properties of the proposed point estimators

In this section we introduce a simulation exercise with the aim of investigating the frequentist properties of point estimators (i.e. posterior means) introduced in section 2 and applied to the analysis of SME survey in section 3. We study bias, mean square error and frequentist coverage of posterior probability intervals.

The simulation is design-based and we do not assume any parametric distribution when generating the data. We create a synthetic population merging the samples of the 2007 and 2008 SME survey discussed in previous section. We drop from the 2007 wave those firms sampled also in the 2008. We obtain a population of size $N = 30451$. Domains are defined by cross-classifying the population by firm size and industry sector; with respect to the data analysis of previous section we collapse the regions. By reducing the number of domains we create sub-populations large enough to be sampled using reasonable sampling rates. Collapsing by region has a milder impact on sub-population skewness with respect to firm size or industry sector. Thereby our synthetic population is divided into $m = 88$ domains whose sizes N_d ($d = 1, \dots, m$) range from 14 to 1339 with an average of 346. The same target parameters (total VA) and auxiliary variable (turnover) studied in section 4 are considered.

We keep this synthetic population as fixed and we repeatedly we repeatedly draw stratified samples with proportional allocation and a 4% sampling rate. Resulting domain-specific sample sizes n_d are adjusted so that $\min(n_d) = 3$. The resulting average domain-specific sample sizes is 14.05, with a maximum of 54. The mean sample size is very close to that of the application.

The Monte Carlo exercise is based on $S = 2000$ repeated samples. Direct estimates and their variances are calculated using analytic formulas. The point estimators we compare in the simulation are:

- i) the posterior mean associated to the Fay-Herriot model specified on the un-transformed scale (UFH). The untransformed Fay-Herriot model can be described as follows: $\hat{Y}_d \sim N(\theta_d, V_d)$, where V_d is the variance of the direct estimator \hat{Y}_d , $\theta_d \sim N(\beta_0 + \beta_1 x_d, \sigma_v^2)$, $\sigma_v \sim Unif(0, A)$, $A = 1000$, $\beta_i \sim N(0, 10^5)$ $i = 0, 1$. We consider the UFH model as it is probably the most “basic” Bayesian model a practitioner would think of for analyzing these data.
- ii) the predictor proposed by Slud and Maiti (2006), we denote as SM;
- iii) the posterior mean obtained from the log-normal model [1]-[3] endowed with the global shrinkage prior [14]. We denote the estimator as LNGS (“Log-Normal with Global Shrinkage”)

1
2
3 1 *iv*)the posterior mean associated to the model [1]-[3] but endowed with local shrinkage prior
4 described in [8]-[10] and $a = 0.5$; we label the estimator as LNLS (“Log-Normal with Local
5 Shrinkage”).
6
7

8 In all cases we set $\beta_0 \sim N(0, 10^5)$, $\beta_1 \sim N(0, 10^5)$, $b_0 = 2$, $c_0 = 1$, as in the application section.
9

10 Let’s denote with est_{ds} the generic estimator calculated for the domain d in the replication s . We
11 compare alternative estimators in terms of relative bias, relative root mean square error and
12 frequentist coverage of probability intervals based on the posterior distribution. Specifically we
13 consider the frequentist coverage defined by the $\alpha/2$ and $(1-\alpha/2)$ quantiles posterior distribution
14 of the target parameter (w.r.t. the coverage probability $1-\alpha$) and set $\alpha = 0.05$. Comparison tools are
15 defined as follows:
16
17
18
19

20
21
22
23 11
$$RBIAS = \frac{1}{m} \sum_{d=1}^m \left| \frac{\frac{1}{S} \sum_{s=1}^S est_{ds} - Y_d}{Y_d} \right|,$$

24
25
26
27 12
$$RRMSE = \frac{1}{m} \sum_{d=1}^m \sqrt{\frac{\frac{1}{S} \sum_{s=1}^S (est_{ds} - Y_d)^2}{Y_d^2}},$$

28
29
30
31 13
$$COV95 = \frac{1}{m} \sum_{d=1}^m \frac{1}{S} \sum_{s=1}^S 1\{Y_d \in (p_{\theta_{ab}|data}(0.025), p_{\theta_{ab}|data}(0.975))\}$$

32
33

34 As we are interested in the frequentist coverage of Bayes estimators, $COV95$ is calculated
35 only for *i*, *iii*, *iv*. Results are summarized in table 4.
36
37
38
39

40 Table 4: Comparison of alternative predictors based on the Monte Carlo experiment

41
42

est	RBIAS	RRMSE	COV95
UFH	0.1573	0.2026	0.708
SM	0.1263	0.1733	--
LNGS	0.1240	0.1719	0.929
LNLS	0.1113	0.1723	0.941

43
44
45
46
47
48
49

50
51
52 20 Results from table 4 show how the posterior means based on the log-normal model with
53 either local or global shrinkage priors and the predictor of Slud and Maiti (2006) perform very
54 closely in terms of mean square error. In terms of bias LNLS is better, but its variance is somewhat
55 bigger, as we can expect from a more flexible, richly parametrized model. Actually, they are based
56 on similar ideas and models, only the priors or the way hyper-parameters are dealt with are
57
58
59
60

1 different, so the results are in line with expectation. We did not expect the bias to be close to 0: in
 2 small area estimation you compromise between the efficiency of a biased synthetic predictor and
 3 the unbiasedness of large variance direct estimators; to some extent estimators associated to areas
 4 with very small area-specific sample sizes are naturally biased.

5 The naïve Fay-Herriot model, specified on the un-transformed scale, performs worse in
 6 terms of both bias and mean square error; the frequentist coverage of the posterior intervals is well
 7 below the 0.95 nominal level. This relatively poor performances reflect the misspecification of the
 8 model, based on the assumption of normality of the direct estimators. It also assumes a linear
 9 relationship between direct estimators and the auxiliary variable on the original scale of the data
 10 (instead of a linear relationship on the log scale); we already noted that this assumption is not
 11 completely unrealistic, so misspecification of the sampling model can be held as responsible for the
 12 not completely satisfying results.

13 The two hierarchical models lead to close performances also in terms of frequentist
 14 coverage of posterior intervals. The advantage of using local shrinkage priors can be appreciated if
 15 we consider the performance for outlying areas, that is those characterized by a deviation from the
 16 synthetic component much larger than most of the remaining areas. We investigate performances
 17 separately for the areas characterized by the larger (on average) model residual. Results related to
 18 the “worst case” area are presented in table 5. We note that this area-specific sample is $n_d = 34$, well
 19 above the average sample size of the simulation.

20
 21 Table 5: Comparison of alternative predictors for area 3, characterized by the largest model
 22 residual

est	<i>RBIAS</i>	<i>RRMSE</i>	<i>COV95</i>
SM	0.4763	0.4790	--
LNGS	0.4414	0.4453	0.505
LNLS	0.0144	0.2515	0.947

23
 24 The LNGS and the SM predictor are based on a similar global shrinkage idea. Results in
 25 table 5 show how the common variance parameter assumed for the random effects cannot
 26 accommodate the “outlier”; the associated model based estimators are severely shrunken towards
 27 the synthetic component: this implies large bias and poor frequentist coverage of the posterior
 28 intervals. Local shrinkage prior associated to LNLS is more flexible and leads to an almost unbiased
 29 predictor and good coverage.

5. Conclusions

We introduced a Bayesian methodology that is useful for small area estimation of means and totals of variables that are positively skewed. This type of variable is often encountered in business surveys. We devote special attention to the specification of a prior distribution for the random effects; our proposal, based on the idea of local shrinkage, is well suited when auxiliary variables with strong predictive power is available, a feature often displayed in business survey data.

The proposed methodology can be easily implemented using widely available MCMC software. Openbugs codes, as well as formulas for the full conditional distributions needed for an independent implementation of the algorithm, are available upon request from the authors.

In summary, we showed that the predictor based on local shrinkage prior has overall acceptable frequentist properties, comparable to the alternatives we consider in the exercise. If most of the areas are well fit by the assumed model and only a minority are outlying, characterized by larger model residuals, we have that, for these areas, local shrinkage priors can lead to estimators with reduced bias and thereby more efficient.

The strategy we propose may be applied to estimating business totals based on any positively skewed variables: value added, turnover, labor cost, and income from sales and services as well as the components of these main aggregates. We discuss the proposed model with reference to real survey data and, more specifically, to the estimation of the total value added, giving consideration to the fact that the value added is the input for calculating important economic aggregates and performance indicators. We address the sub-populations of Italian small- and medium-sized manufacturing firms classified according to sub-national region, industry and firm-size classes. We limit our attention to small and medium enterprises, that is, on firms with less than 100 employees because in general, as well as in Italy, larger firms are censused, and small area estimation is therefore not needed.

This research can be extended and complemented in many directions: One important problem not considered here is that of benchmarking of small area estimates to known totals for more aggregate domains. A second aspect to address is the longitudinal extension of the model specification in order to borrow strength not only from covariates but also information repeated over time. This also makes it possible to produce estimates at different time points.

References

- Barndorff-Nielsen, O.E. (1977). Exponentially decreasing distributions for the logarithm of particle size, *Proceedings of the Royal Statistical Society, series A*, 353, 401-419.
- Breinlich, H., Ottaviano, G.I.P., & Temple, J.R.W. (2014). Regional Growth and Regional Decline. in: *Handbook of Economic Growth*, edition 1, vol. 2, ch. 4, 683-779 Elsevier.
- Burgard, J.P., Munnich R., & Zimmermann T. (2014). The Impact of Sampling Designs on Small Area Estimates for Business Data. *Journal of Official Statistics*, 30, 4, 749-771.
- Carlin, B.P., & Louis, T.A. (2000). *Bayes and empirical Bayes data analysis*. New York, Chapman and Hall.
- Cobb B.R., Rumi R. & Salmeron A. (2012). Approximating the distribution of a sum of log-normal random variables. In: *Proc. of the VI European Workshop on Probabilistic Graphical Models*, Granada.
- Datta, G.S. & Lahiri, P. (1995). Robust hierarchical Bayes estimation of small area characteristics in presence of covariates and outliers. *Journal of Multivariate Analysis*, 54, 2, 310-328.
- Datta, G. S., Hall, P. and Mandal, A. (2011). Model selection by testing for the presence of small-area effects in area-level data. *Journal of the American Statistical Association*, 106, 362-374.
- Datta, G.S., Mandal A. (2015) Small Area Estimation with Uncertain Random Effects, *Journal of the American Statistical Association*, in press.
- Eurostat (2011), *Key figures on European business - with a special feature on SMEs*, Eurostat Pocketbooks.
- Eurostat (2015a). *Eurostat regional yearbook 2015. Statistical books, General and regional statistics*.
- Eurostat (2015b). *Regions in the European Union - Nomenclature of territorial units for statistics - NUTS 2013/EU-28*.
- Fabrizi, E., Ferrante, M.R., Pacei, S. & Trivisano C. (2011). Hierarchical Bayes multivariate estimation of poverty rates based on increasing thresholds for small domains. *Computational Statistics & Data Analysis*, 55, 1736 – 1747.
- Fabrizi, E. & Trivisano C. (2016). Small area estimation of the Gini concentration coefficient. *Computational Statistics & Data Analysis*, 99, 223 – 234.
- Fay, R. & Herriot, R. (1979). Estimates of income for small places: an application of James-Stein procedures to census data. In: *Journal of the American Statistical Association*, 74, p. 269-277.
- Fenton, L. (1960) The Sum of Log-Normal Probability Distributions in Scatter Transmission Systems. *IRE Transactions on Communication Systems* 8(1), 57-67.

1
2
3 1 Ferrante, M.R., & Trivisano C. (2010). Small area estimation of the number of firms'
4 recruits by using multivariate models for count data. *Survey Methodology*, 36, 2, 171-180.

5
6 3 Fruhwirth-Schnatter, S., & Wagner, H. (2010). Bayesian variable selection for random
7 intercept modelling of Gaussian and non-Gaussian data. In J. Bernardo, M. Bayarri, J.O. Berger,
8 A.P. Dawid, D. Heckerman, A.F.M. Smith, M. West (eds.), *Bayesian Statistics*, 9, 165–200, Oxford
9 Univ. Press.

10
11 7 Gelman A., Meng X., Stern H. (1996) Posterior predictive assessment of model fitness via
12 realized discrepancies. *Statistica Sinica*, 6, 733-807.

13
14 9 Gelman A. (2006). Prior distributions for variance parameters in hierarchical models.
15 *Bayesian Analysis*, 1, 3, 515–533.

16
17 11 Griffin, J.E., & Brown, P.J. (2010). Inference with normal-gamma prior distributions in
18 regression problems. *Bayesian Analysis*, 5, 171- 188.

19
20 13 Horrace W.C. (2005) Some results on the multivariate truncated normal distribution, *Journal*
21 *of Multivariate Statistics*, 94, 209-221.

22
23 15 Ibrahim J., Chen M., Sinha D. (2001). *Bayesian Survival Analysis*. Springer-Verlag, New
24 York.

25
26 17 Jiang J. & Lahiri P. (2005) Mixed model prediction and small area estimation (with
27 discussion), *TEST*, 15, 1-96

28
29 19 Karlberg F. (2000) Survey estimation for highly skewed populations in the presence of
30 zeroes. *Journal of Official Statistics*, 16, 229-241.

31
32 21 Mazmanyany, L., Ohanyan, V., Trietsch, D. (2009) The lognormal central limit theorem for
33 positive random variables. Reproduced as Appendix in Baker K.R, Trietsch D. (2009), *Principles of*
34 *sequencing and scheduling*, John Wiley and Sons, New York.

35
36 24 Militino, A.F., Ugarte, M.D. & Goicoa T. (2015). Deriving small area estimates from
37 information technology business surveys. *Journal of the Royal Statistical Society A*, 178, 1051-
38 1067.

39
40 27 Park, Casella G. (2008). The Bayesian lasso, *Journal of the American Statistical*
41 *Association*, 103, 681–686.

42
43 29 Pratesi M. ed. (2016). *Analysis of Poverty Data by Small Area Estimation*, John Wiley, New
44 York.

45
46 31 OECD (2013), *OECD Regions at a Glance 2013*: OECD Publishing.
47 http://dx.doi.org/10.1787/reg_glance-2013-en

48
49 33 Pfeffermann D. (2014). Small Area Estimation. In *International Encyclopedia of Statistical*
50 *Science*, M. Lovric (eds.), 1346-1349, Springer-Verlag.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1 Rao J.N.K. (2003). Small area estimation, John Wiley and Sons. New York.

2 Rivière P. (2002): What Makes Business Statistics Special? *International Statistical Review*,

3 70, 1, 145-159.

4 Slud, E. V. & Maiti, T. (2006). Mean-squared error estimation in transformed Fay–Herriot

5 models. *Journal of the Royal Statistical Society, Series B*, 68, 2, 239–257.

6 Statistics Canada (2007) 2005 Survey of Financial Security. Public Use Microdata File User

7 Guide. Ottawa: Statistics Canada. (Available from

8 <http://www.statcan.gc.ca/pub/13f0026m/13f0026m2007001-eng.htm>.)

9 Wolter K.M. (1985). Introduction to variance estimator, New York, Springer-Verlag.

10 You, Y. & Rao, J.N.K. (2002). Small area estimation using unmatched sampling and linking

11 models. *Canadian Journal of Statistics*, 30, 3–15.

12

1
2
3 **Appendix 1**

4
5
6 **Proof of [6]**

7 To start with we note that

8
9
10
$$\hat{\theta}_d^{B1} = \exp\left\{\mathbf{x}'_d \boldsymbol{\beta} + \gamma_d (\hat{Z}_d - \mathbf{x}'_d \boldsymbol{\beta}) + \frac{1}{2} \gamma_d \hat{\delta}_d\right\} = \exp\left\{\gamma_d \hat{Z}_d + (1 - \gamma_d) \mathbf{x}'_d \boldsymbol{\beta} + \frac{1}{2} \gamma_d \hat{\delta}_d\right\}$$
 so

11
12
13
$$E\left[\exp\left\{\gamma_d \hat{Z}_d + (1 - \gamma_d) \mathbf{x}'_d \boldsymbol{\beta} + \frac{1}{2} \gamma_d \hat{\delta}_d\right\}\right] = \exp\left\{(1 - \gamma_d) \mathbf{x}'_d \boldsymbol{\beta} + \frac{1}{2} \gamma_d \hat{\delta}_d\right\} E\left[\exp\{\gamma_d \hat{Z}_d\}\right].$$

14
15
16 We note that $E(\hat{Z}_d) = E_M\{E(\hat{Z}_d | \eta_d)\} = E_M\{\eta_d\} = \mathbf{x}'_d \boldsymbol{\beta}$; analogously

17
18
19
$$V(\hat{Z}_d) = V_M\{E(\hat{Z}_d | \eta_d)\} + E_M\{V(\hat{Z}_d | \eta_d)\} = \psi_d + \hat{\delta}_d$$
 . As a consequence $E(\gamma_d \hat{Z}_d) = \gamma_d \mathbf{x}'_d \boldsymbol{\beta}$,

20
21
22
$$V(\gamma_d \hat{Z}_d) = \gamma_d^2 (\psi_d + \hat{\delta}_d) = \psi_d \gamma_d$$
 and $E\left[\exp\{\gamma_d \hat{Z}_d\}\right] = \exp\left\{\gamma_d \mathbf{x}'_d \boldsymbol{\beta} + \frac{1}{2} \psi_d \gamma_d\right\}$. This leads to

23
24
25
$$E(\hat{\theta}_d^{B1}) = \exp\left\{\mathbf{x}'_d \boldsymbol{\beta} + \frac{1}{2} \gamma_d (\psi_d + \hat{\delta}_d)\right\} = \exp\left\{\mathbf{x}'_d \boldsymbol{\beta} + \frac{1}{2} \psi_d\right\}$$
 that coincides with $E_M(\theta_d)$.

26
27
28
29
30
31

32 **Proof of [7]**

33 Let's introduce some matrix notation. Let $\mathbf{z} = \text{vec}(\hat{Z}_d)$ is the vector containing, $\boldsymbol{\Psi} = \text{diag}(\psi_d)$,

34 $\boldsymbol{\Delta} = \text{diag}(\hat{\delta}_d)$ the diagonal matrixes containing the variance components; let \mathbf{X} be the matrix with

35
36
37
38
39 rows \mathbf{x}'_d , $d = 1, \dots, m$.

40
41
42 Standard Bayesian analysis of normal linear mixed model lead to $\boldsymbol{\beta} | \mathbf{z}, \boldsymbol{\Psi} \sim N(\hat{\boldsymbol{\beta}}_{gls}, \mathbf{V}(\boldsymbol{\Psi}))$ where

43
44
45
$$\hat{\boldsymbol{\beta}}_{gls} = [\mathbf{X}' (\boldsymbol{\Psi} + \boldsymbol{\Delta})^{-1} \mathbf{X}]^{-1} \mathbf{X}' (\boldsymbol{\Psi} + \boldsymbol{\Delta})^{-1} \mathbf{z}$$
 , $\mathbf{V}(\boldsymbol{\Psi}) = [\mathbf{X}' (\boldsymbol{\Psi} + \boldsymbol{\Delta})^{-1} \mathbf{X}]^{-1}$. We can calculate

46
47
48
$$E(\boldsymbol{\eta} | \mathbf{z}, \boldsymbol{\Psi}) = E_{\boldsymbol{\beta} | \mathbf{z}, \boldsymbol{\Psi}} E(\boldsymbol{\eta} | \mathbf{z}, \boldsymbol{\beta}, \boldsymbol{\Psi})$$
 where $\boldsymbol{\eta} = \text{vec}(\eta_d)$. $E_{\boldsymbol{\beta} | \mathbf{z}, \boldsymbol{\Psi}} E(\boldsymbol{\eta} | \mathbf{z}, \boldsymbol{\beta}, \boldsymbol{\Psi}) = E_{\boldsymbol{\beta} | \mathbf{z}, \boldsymbol{\Psi}} (\mathbf{X} \boldsymbol{\beta}) = \mathbf{X} \hat{\boldsymbol{\beta}}_{gls}$.

49
50
51 Analogously $V(\boldsymbol{\eta} | \mathbf{z}, \boldsymbol{\Psi}) = V_{\boldsymbol{\beta} | \mathbf{z}, \boldsymbol{\Psi}} E(\boldsymbol{\eta} | \mathbf{z}, \boldsymbol{\beta}, \boldsymbol{\Psi}) + E_{\boldsymbol{\beta} | \mathbf{z}, \boldsymbol{\Psi}} V(\boldsymbol{\eta} | \mathbf{z}, \boldsymbol{\beta}, \boldsymbol{\Psi})$. If we denote the vector of small

52
53
54 area predictors (on the log scale) conditional on $\boldsymbol{\beta}$ as $\boldsymbol{\eta}^{B1} = E(\boldsymbol{\eta} | \mathbf{z}, \boldsymbol{\beta}, \boldsymbol{\Psi}) = \boldsymbol{\Gamma} \mathbf{z} + (\mathbf{I} - \boldsymbol{\Gamma}) \mathbf{X} \boldsymbol{\beta}$ with

55
56
57
$$\boldsymbol{\Gamma} = \boldsymbol{\Psi} (\boldsymbol{\Psi} + \boldsymbol{\Delta})^{-1}$$
 we have that $V(\boldsymbol{\eta} | \mathbf{z}, \boldsymbol{\Psi}) = V_{\boldsymbol{\beta} | \mathbf{z}, \boldsymbol{\Psi}} \boldsymbol{\Gamma} \mathbf{z} + (\mathbf{I} - \boldsymbol{\Gamma}) \mathbf{X} \boldsymbol{\beta} + E_{\boldsymbol{\beta} | \mathbf{z}, \boldsymbol{\Psi}} \mathbf{G}_1$ with $\mathbf{G}_1 = \boldsymbol{\Psi} (\boldsymbol{\Psi} + \boldsymbol{\Delta})^{-1} \boldsymbol{\Delta}$

58
59
60 . Taking expectation with respect to $p(\boldsymbol{\beta} | \mathbf{z}, \boldsymbol{\Psi})$ we get

1
2
3 1 $V(\boldsymbol{\eta} | \mathbf{z}, \boldsymbol{\Psi}) = (\mathbf{I} - \boldsymbol{\Gamma}) \mathbf{X}' \boldsymbol{\beta} \mathbf{V}(\boldsymbol{\Psi}) \mathbf{X} (\mathbf{I} - \boldsymbol{\Gamma}) + \mathbf{G}_1 = \mathbf{G}_2 + \mathbf{G}_1$. We note that $p(\boldsymbol{\eta} | \mathbf{z}, \boldsymbol{\Psi})$ is a multivariate
4
5 2 normal distribution. If we consider an individual η_d have that $E(\eta_{ijr} | \mathbf{z}, \boldsymbol{\Psi}) = \mathbf{x}_d' \hat{\boldsymbol{\beta}}_{gls}$,

6
7
8 3
$$V(\eta_d | \mathbf{z}, \boldsymbol{\Psi}) = \gamma_d \hat{\delta}_d + (1 - \gamma_d)^2 \mathbf{x}_d' \left(\sum_{d=1}^m \frac{1}{\psi_d + \hat{\delta}_d} \mathbf{x}_d \mathbf{x}_d' \right)^{-1} \mathbf{x}_d = g_{1,d} + g_{2,d}.$$

9
10
11 4 As $\theta_d = \exp(\eta_d)$ formula [7] follows.

12 5

13 6

14
15
16
17 7 **Proof of [11] and subsequent statements**

18
19 8 From [9] we have that $p(\boldsymbol{\Psi} | \lambda) = \prod_{d=1}^m \frac{\lambda^a}{\Gamma(a)} \psi_d^{a-1} e^{-\lambda \sum_{d=1}^m \psi_d} \propto \lambda^{ma} e^{-\lambda \sum_{d=1}^m \psi_d} \prod_{d=1}^m \psi_d^{a-1}$

20
21
22 9 Conditioning on a is omitted as it is treated as a known constant; m is a shortcut notation for the
23
24 10 overall number of areas.

25
26 11 We can obtain the marginal prior $p(\boldsymbol{\Psi})$ using the integral $p(\boldsymbol{\Psi}) = \int_{\mathbb{R}^+} p(\boldsymbol{\Psi} | \lambda) p(\lambda) d\lambda$. As

27
28
29 12 $\lambda \sim \text{Gamma}(b_0, c_0)$ we straightforwardly get

30
31 13
$$p(\boldsymbol{\Psi}) \propto \frac{\Gamma(ma + b_0)}{\left(\sum_{d=1}^m \psi_d + c_0 \right)^{ma + b_0}} \prod_{d=1}^m \psi_d^{a-1}$$

32
33
34
35 14 Applying the transformation $\tau_d = \sqrt{\psi_d}$ on each component of $\boldsymbol{\Psi}$ we obtain

36
37
38 15
$$p(\boldsymbol{\tau}) \propto \left(\sum_{d=1}^m \tau_{ijr}^2 + c_0 \right)^{-(ma + b_0)} \prod_{d=1}^m \tau_d^{2a-1}$$

39
40
41
42 16 For the special case $a = \frac{1}{2}$ the density of $p(\boldsymbol{\tau})$ simplifies to $p(\boldsymbol{\tau}) \propto \left(\sum_{d=1}^m \tau_d^2 + c_0 \right)^{-(ma + b_0)}$ or

43
44
45 17 equivalently to

46
47 18
$$p(\boldsymbol{\tau}) \propto \left(1 + \frac{1}{c_0} \sum_{d=1}^m \tau_d^2 \right)^{-\left(\frac{m}{2} + b_0\right)},$$

48
49
50
51 19 $\tau_d > 0, \forall d$. This expression can be recognized as the kernel of the density of a multivariate half-t
52
53 20 distribution with mean vector 0 and diagonal scale matrix. A multivariate half-t is a multivariate t
54
55 21 for which we apply the absolute value transformation on each component. We can then write
56
57 22 formula [13].

58
59
60 23

1 To prove that univariate priors $p(\tau_d)$ let's start from $p(\tau_d) = \int \dots \int p(\boldsymbol{\tau}) d\boldsymbol{\tau}_{-d}$. We can represent as
 2 $p(\boldsymbol{\tau})$ as the result of applying the absolute value transformation on a multivariate t distribution, i.e

3 $\boldsymbol{\tau} = |\boldsymbol{\tau}^*|$ with $p(\boldsymbol{\tau}^*) = \int p(\boldsymbol{\tau}^* | \xi) p(\xi) d\xi$ where $p(\boldsymbol{\tau}^* | \xi) \sim MVN\left(\mathbf{0}, \frac{2b_0}{c_0} \xi \mathbf{I}_m\right)$ and

4 $p(\xi) \sim IGamma(b_0, b_0)$. We use the fact that a random vector distributed according to a
 5 multivariate t distribution can be represented as an inverse Gamma mixture of a multivariate
 6 normal.

7 As the variance covariance matrix of $\boldsymbol{\tau}^*$ is diagonal $p(\boldsymbol{\tau}^*) = \int \prod_{d=1}^m p(\tau_d^* | \xi) p(\xi) d\xi$.

8 Horrace (2005) studies truncated multivariate normal distributions and shows that univariate
 9 marginal distributions from a multivariate half-Normal (obtained applying the absolute value
 10 transformation on each component) are univariate half-Normals if and only the variance-covariance
 11 matrix of the parent multivariate normal is diagonal. As a consequence

$$12 \quad p(\tau_d) = \int \dots \int \left[\int \prod_{d=1}^m p(\tau_d | \xi) p(\xi) d\xi \right] d\boldsymbol{\tau}_{-d}$$

13 where each $p(\tau_d | \xi)$ is distributed as an half-Normal distribution.

14 If we change the order of integration and use conditional independence of $p(\tau_d | \xi)$ we obtain that

15 $p(\tau_d)$ are marginally half-t distributed.

16 To prove that τ_d are linearly independent of each other we write

$$17 \quad V(\boldsymbol{\tau}) = E_{\xi} \{V(\boldsymbol{\tau} | \xi)\} + V_{\xi} \{E(\boldsymbol{\tau} | \xi)\}$$

18 and note that $E(\boldsymbol{\tau} | \xi) = \mathbf{0}$ while $E_{\xi} \{V(\boldsymbol{\tau} | \xi)\} = E(\xi) \frac{2b_0}{c_0} \mathbf{I}_m$, which is of course diagonal.