

# Small area solutions for the analysis of pollution data

Daniela Cocchi, Enrico Fabrizi, Carlo Trivisano

Università di Bologna

TIES 2002

Genova, June 22nd 2002

## Summary

The small area problem in the environmental context

Theoretical framework

The linear hierarchical model

Special case: ANOVA model

Approximating criterion: least squares

Prior conjectures

Statistic chosen for conditioning

The normal case and the departure from normality

Comments on least squares approximations

A simulated experiment

Consequences of misspecifications in prior evaluations

High applicability to environmental problems

Extensions

## The small area problem in the environmental context

Simultaneous estimation of parameters related to different subpopulations (domains) of a more general population

Desired property: additivity

### Example 1

Erosion from agricultural land in a watershed (Opsomer, Botts, Kim, 2001)

Population of areas (160 acres plot)

Auxiliary variables available

Mixed effect linear model

## Example 2

### Emissions inventories

(CORINAIR project in EU)

national estimates of emission volumes (by pollutant)

many domain estimates

spatial

economic sectors

time evolution

Building inventories

**Census** of major pollution sources

**Indirect estimation** of small and very diffuse sources

Activity indicators

Emission factors

## Hierarchical classification of emission activities SNAP system (Selected Nomenclature for Air Pollution)

$$E_{p,a,t} = A_{a,t} F_{p,a}$$

$A_{a,t}$ : volume of activity  $a$  in period  $t$  subject to sampling random variation;  $F_{p,a}$ : emission factor for pollutant  $p$  in activity  $a$ ; seldom estimated on a sample basis.

$$E_{p,t} = \sum_{a=1}^N E_{a,p,t}$$

**Approach: bottom up.**

Suggestion: ANOVA model is a realistic proposal for  $A_{a,t}$  since it does not use auxiliary information which is very heterogeneous and suffers from problems of spatio-temporal definition.

## Theoretical framework

Typical of the finite population context

It can be solved within

The design based approach

parameters: unknown characteristics of the subpopulations

The model based approach

where hyperparameters are considered

We discuss:

Model based approach

Hierarchical modeling

    stressing Bayesian solutions

Two main motivations

A natural way of building and solving models

Managing different aggregation level of the information

## The linear hierarchical model

It is a **structural model** for the population

$$\begin{aligned} \eta &= Z_2 \beta + \varepsilon & \eta, \varepsilon &: N \times 1 & Z_2 &: N \times p & \beta &: p \times 1 \\ \beta &= Z_1 \beta_0 + \delta & Z_1 &: p \times q & \beta_0 &: q \times 1 & \delta &: q \times 1 \end{aligned}$$

where

$$Z_2 = \text{diag}[i_{Nk}]$$

is the **small area indicator matrix** and the **incidental parameter**  $\beta$  is the vector of **small area means** object of inference.



Assumptions on the model (only on the first two moments of the distributions):

$$C(\delta', \varepsilon | Z_1, Z_2) = 0$$

$$E(\varepsilon | Z_2, \sigma_\varepsilon^2, \beta) = E(\varepsilon | \sigma_\varepsilon^2) = 0$$

$$V(\varepsilon | Z_2, \sigma_\varepsilon^2, \beta) = V(\varepsilon | \sigma_\varepsilon^2) = \sigma_\varepsilon^2 V$$

$$E(\delta | Z_1, \sigma_\delta^2, \beta_0) = E(\delta | \sigma_\delta^2) = 0$$

$$V(\delta | Z_1, \sigma_\delta^2, \beta_0) = V(\delta | \sigma_\delta^2) = \sigma_\delta^2 B_0$$

A set of assumptions of conditional independence is able to simplify the model.

Many small area models are special cases of this general model.

## Role of covariates

a) At the individual level

The slightly different **Battese et. al. (1988)** model

$$\eta = X\beta + v$$

$$v = Z_2\delta + \varepsilon$$

$$V = I_N \quad B_0 = i_p i_p' \quad Z_1 = I_p$$

with a  $\beta$  value on the whole population and small area random effects

b) At the area level

$$N = p$$

$$\eta, \varepsilon : p \times 1 \quad Z_2 = I_p \quad \beta : p \times 1$$

**Fay and Herriott (1979) model** adds distributional assumptions

$$e_k \sim N(0, D_k), \quad k = 1, \dots, p$$

$$\delta_k \sim N(0, A), \quad k = 1, \dots, p$$

Frequentist solutions

Empirical Best Linear Unbiased Prediction.

Estimating  $\sigma_\delta^2$  and  $\sigma_\varepsilon^2$ , and introducing the estimates in the Best Linear Unbiased Predictor of  $\beta$ .

## Special case: ANOVA model

Some basic simplifications

$$V = I_N$$

$$Z_2 Z_1 = i_N$$

## A comprehensive statistical model

A **sampling model** can be added to the **structural model**

$$S = (e'_{s_1}, \dots, e'_{s_n})' : (n \times N)$$

where  $e_{s_i}$  is the  $s_i$  column of  $I_N$ .

## Approximating criterion: least squares

### The general least squares principle

Approximation to a normal distribution with the same first two moments of the exact one.

Finite populations: this approximation is conditional on  $(Z, S)$ .

Normal approximation on  $(\beta, t|Z, S)$  and not on the whole  $(\beta, t, Z, S)$ .

TYPE OF INFERENCE: Bayesian

Why approximated solutions when MCMC solutions which approximated posterior distributions are easily available?

Because they are analytically approximated and need the elicitation of a relatively small number of prior guesses.

Be  $X$  the data and  $\psi$  the parameter, for any  $\beta = g(\psi)$  and  $t = t(X)$ , we define  $l(t) = \{a't\}$ .

If:

$$l_{Z,S}(t) = \{a(z,s)'t\}$$

$$\begin{aligned} E_{LS}^{Z,S}(\beta|t) &= \arg \min_l E\left(\|\beta - l_{Z,S}(t)\|^2 | Z, S\right) \\ &= \arg \min_l E\left(\|E(\beta|t, Z, S) - l_{Z,S}(t)\|^2 | Z, S\right) \\ &= E(\beta|Z, S) + C(\beta, t'|Z, S)\{V(t|Z, S)\}^{-1}\{t - E(t|Z, S)\} \end{aligned}$$

$$\begin{aligned} V_{LS}^{Z,S}(\beta|t) &= \min_{l_{Z,S}} E\left(\|E(\beta|t, Z, S) - l_{Z,S}(t)\|^2 | Z, S\right) \\ &= V(\beta|Z, S) + C(\beta, t'|Z, S)\{V(t|Z, S)\}^{-1}C(t, \beta'|Z, S) \end{aligned}$$

(Cocchi and Mouchart, 1996)

## Prior conjectures

Only on a certain number of moments

Conjecturing on moments is more immediate than conjecturing on distributional assumptions

Their type and number depend on the statistic on which conditioning

## Statistic chosen for conditioning

Its choice can enrich the solution and determines its complexity

a) **Solution conditional on**

$t = T_0 = Z_2' S' y$ : vector of small area totals

Priors to be elicited

$$E(\beta_0 | Z_1) = E(\beta_0) = b_0$$

$$V(\beta_0 | Z_1) = V(\beta_0) = M_0$$

$$v_\delta = E(\sigma_\delta^2 | Z_1) = E(\sigma_\delta^2)$$

$$v_\varepsilon = E(\sigma_\varepsilon^2 | Z_2) = E(\sigma_\varepsilon^2)$$



## b) Solution conditional on

$$t = T = (T_0, T_1, T_2)$$

i.e. the vector of small area totals and the sum of squares

$T_\delta$ : between small areas and

$T_\varepsilon$ : within small areas

Since  $T$  contains a polynomial of order  $d$ , prior information on moments of order  $2d$  must be used.

For computing the solution the following moments are needed

$$\alpha_j = E\left(\delta_k^j \mid Z_1, \beta_0, \sigma_\delta^2\right), \quad j=3,4$$

$$\varphi_j = E\left(\varepsilon_k^j \mid Z_2, \beta, \sigma_\varepsilon^2\right), \quad j=3,4$$

New priors to be elicited besides the group of priors above:

$$V_\delta = V(\sigma_\delta^2 | Z, S) \quad V_\varepsilon = V(\sigma_\varepsilon^2 | Z, S)$$

$$a_3 = E(\alpha_3 | Z, S) \quad a_4 = E(\alpha_4 | Z, S)$$

$$f_3 = E(\varphi_3 | Z, S) \quad f_4 = E(\varphi_4 | Z, S)$$

$$c_{0,1} = C(\beta_0, \sigma_\delta^2 | Z, S)$$

It is not a normal approximation!

It does not mean approximating normal distributions with the same first two moments: **the model may contain moments up to the 4<sup>th</sup>.**

Gain: this solution can consider **asymmetry** and **kurtosis**.

## Solution

a) For  $t=T_0$

$$\begin{aligned} E_{LS}^{Z,S}(\beta|T_0) &= \mathbf{a} \left[ gb_0 + (1-g)\bar{\bar{y}} \right] + (I_p - \Delta_{\mathbf{a}}) \mathbf{y} \\ &= \mathbf{w}v_{\varepsilon} \left[ gb_0 + (1-g)\bar{\bar{y}} \right] + v_{\varepsilon} (v_{\varepsilon}I_p - \Delta_{\mathbf{w}}) \mathbf{y} : \quad p \times 1 \end{aligned}$$

$$\begin{aligned} V_{LS}^{Z,S}(\beta|T_0) &= v_{\delta} \Delta_{\mathbf{a}} + M_0 \mathbf{g} \mathbf{a} \mathbf{a}' \\ &= \mathbf{w}v_{\varepsilon} \Delta_{\mathbf{w}} + M_0 v_{\varepsilon}^2 \mathbf{g} \mathbf{w} \mathbf{w}' : \quad p \times p \end{aligned}$$

where

$$\mathbf{a} = [a_k] = \left[ v_{\delta}^{-1} \left( v_{\delta}^{-1} + n_k v_{\varepsilon}^{-1} \right)^{-1} \right] = \left[ v_{\varepsilon} \left( v_{\varepsilon} + n_k v_{\delta} \right)^{-1} \right] : \quad p \times 1$$

$$\mathbf{w} = [w_k] = v_\varepsilon^{-1} [a_k] = \left[ (v_\varepsilon + n_k v_\delta)^{-1} \right] : p \times 1$$

$$\Delta_{\mathbf{a}} = \text{diag}[a_k] : p \times p \quad \Delta_{\mathbf{w}} = \text{diag}[w_k] = (v_\varepsilon + n_k v_\delta)^{-1} : p \times p$$

$$\bar{y} = \frac{\sum_k \frac{n_k \bar{y}_k}{v_\varepsilon + n_k v_\delta}}{\sum_k \frac{n_k}{v_\varepsilon + n_k v_\delta}} = \frac{v_\varepsilon \mathbf{w}' \Delta_{\mathbf{n}} \mathbf{y}}{v_\varepsilon \mathbf{w}' \Delta_{\mathbf{n}} \mathbf{i}_p} = \frac{\mathbf{w}' T_0}{\mathbf{w}' \mathbf{n}}$$

$$g = (1 + M_0 \mathbf{w}' \mathbf{n})^{-1} = M_0^{-1} (M_0^{-1} + \mathbf{w}' \mathbf{n})^{-1}$$

$$\Delta_{\mathbf{n}} = Z_2' S' S Z_2 = \text{diag}[n_k] : p \times p \quad \mathbf{n} = \Delta_{\mathbf{n}} \mathbf{i}_p : p \times 1 \quad \mathbf{y} = \Delta_{\mathbf{n}}^{-1} T_0 : p \times 1$$

$$y = n^{-1} \mathbf{i}'_n \mathbf{y} = n^{-1} \mathbf{i}'_p T_0 = n^{-1} \mathbf{n}' \mathbf{y}$$

b) For  $t=T$

heavier calculations, no elegant formula because of the difficulty of writing analytically

$$\left[ V(T|Z, S) \right]^{-1} = \begin{bmatrix} \mathbf{V}^{00} & \mathbf{v}^{01} & \mathbf{v}^{02} \\ \mathbf{v}^{10} & \mathbf{v}^{11} & \mathbf{v}^{12} \\ \mathbf{v}^{20} & \mathbf{v}^{21} & \mathbf{v}^{22} \end{bmatrix} : (p+2) \times (p+2)$$

$$E_{LS}^{Z,S}(\beta|T) = (I_p - C_*) i_p b_0 + C_* \mathbf{y} + C_{**} \begin{bmatrix} T_\delta - \left[ v_\varepsilon (p-1) + v_\delta (n - n^{-1} \mathbf{n}' \mathbf{n}) \right] \\ T_\varepsilon - v_\varepsilon (n-p) \end{bmatrix}$$

where

$$C_* = \left[ C_0 V^{00} + c_1 (\mathbf{v}^{01})' \right] \Delta_{\mathbf{n}} : p \times p$$

$$C_{**} = \begin{bmatrix} C_0 \mathbf{v}^{01} + c_1 \mathbf{v}^{11} & C_0 \mathbf{v}^{02} + c_1 \mathbf{v}^{12} \end{bmatrix} : p \times 2$$

after defining:

$$C(\beta|T') = \begin{bmatrix} C_0 & c_1 & 0 \end{bmatrix} : p \times (p+2)$$

## The normal case and the departure from normality

In case of symmetry: 
$$\left[ V(T|Z, S) \right]^{-1} = \begin{bmatrix} \left[ V(T_0|Z, S) \right]^{-1} & 0 \\ 0 & \left[ V \begin{pmatrix} T_\delta \\ T_\varepsilon \end{pmatrix} | Z, S \right]^{-1} \end{bmatrix}$$

when normality holds:  $C[T_\delta, T_\varepsilon] = V_\varepsilon (p-1)(n-p)$

When  $c_{0,1} = C(\beta, T_\delta | Z, S) = 0$

Then  $C_{**} = 0$  and  $C(\beta | T' | Z, S) \left[ V(T|Z, S) \right]^{-1} = C(\beta | T'_0 | Z, S) \left[ V(T_0|Z, S) \right]^{-1}$

As a consequence

$$E_{LS}^{Z,S}(\beta | T) = E_{LS}^{Z,S}(\beta | T_0)$$

## Comments on least squares approximations

They borrow strength from the other subgroups and from the prior conjectures

The idea recalls two important points :

posterior linearity (for estimating the general mean starting from the least squares approximations of the small area parameters)

empirical Bayes solutions.

The solution is robust

since it depends only on moments.

The solution seems not to have practical drawbacks linked to the number of small areas or to the fact of not having evidence for some subgroups.

## A simulated experiment

**Aim:** Checking the improvement of conditioning on  $T$  instead of  $T_0$  in the approximated posterior expectations.

**Data generation process:** random generation of  $\beta$  and  $\eta$  from

normal - normal

lognormal - lognormal.

Lognormal parameters are set in order to have the same first two moments of the normal distribution.

**Exact Bayesian least squares solutions**

i.e. right prior conjectures.



## Design of the experiment

1. generation of 1000 populations of  $N= 1000$  elements with mean  $\beta_0=1$ .

under 4 different hypotheses on structural variability:

$$\sigma_\delta = 1, \sigma_\varepsilon = 0.5 \qquad \sigma_\delta = 0.5, \sigma_\varepsilon = 0.5$$

$$\sigma_\delta = 2, \sigma_\varepsilon = 1 \qquad \sigma_\delta = 1, \sigma_\varepsilon = 1$$

and 2 different hypotheses on the domains:  $p=10$ ;  $p=40$ ;

2. generation of 100 samples from each population with

$$n_k = 10, k=1, \dots, p \text{ when } p=10 \qquad n_k = 5, k=1, \dots, p \text{ when } p=40;$$

3. evaluation of the performances of the posterior expectations by means of mean square errors averaged over the domains.

## Lognormal – Lognormal

|  | $p=10 \ n_k=20$ | $p=40 \ n_k=5$ |
|--|-----------------|----------------|
| $\sigma_\delta=1 \ \sigma_\varepsilon=0.5$   | 1.043           | 1.088          |
| $\sigma_\delta=0.5 \ \sigma_\varepsilon=0.5$ | 1.012           | 1.023          |
| $\sigma_\delta=2 \ \sigma_\varepsilon=1$     | 1.085           | 1.132          |
| $\sigma_\delta=2 \ \sigma_\varepsilon=2$     | 1.031           | 1.054          |

## Normal – Normal ( $c_{0,1} = 1$ )

|  | $p=10 \ n_k=20$ | $p=40 \ n_k=5$ |
|--|-----------------|----------------|
| $\sigma_\delta=1 \ \sigma_\varepsilon=0.5$   | 1.012           | 1.021          |
| $\sigma_\delta=0.5 \ \sigma_\varepsilon=0.5$ | 1.001           | 1.009          |
| $\sigma_\delta=2 \ \sigma_\varepsilon=1$     | 1.023           | 1.044          |
| $\sigma_\delta=2 \ \sigma_\varepsilon=2$     | 1.011           | 1.014          |

## Consequences of misspecifications in prior evaluations

a) Not important: errors in conjectures on variances

$$V(\beta_0) \quad V(\sigma_\delta^2) \quad V(\sigma_\varepsilon^2)$$

b) As expected: errors in conjectures on  $E(\beta_0)$

c) The most dangerous: errors in conjectures on expectations of variances

$$E(\sigma_\delta^2) \quad E(\sigma_\varepsilon^2)$$

d) Exchanging the right lognormal and normal conjectures

For normal-normal data generation: the most dangerous are the errors on the 3rd moment.

For lognormal - lognormal data generation: the most dangerous are the errors on the 4<sup>th</sup> moments.

## High possibility of application to environmental problems

Suitable solution for estimating emission inventories since

- a) the compositeness and heterogeneity of local emissions is so high that a random effect model is in many cases the only practicable solution
- b) all methods for building emissions inventories involve expert evaluations, which can be managed as priors.

## Extensions

introduction of covariates

relationships with the design