

Methodological proposals I: general framework, ecological sampling, and geostatistics

WP5 Mobile Phone Data

E. Esteban, S. Saldaña, **D. Salgado**, L. Sanguiao

Statistics Spain (INE)

Madrid, 7-8 June, 2017



Overview

- **Definition** of the problem
- **Framework**: general lines
- More concrete **proposals**
- Diverse **remarks**

Problem

Framework

Concrete proposals

Remarks



The problem: data sets

- **Two type of data sets:**

- Data at **mobile device** level (CDR/DDR, signaling, ...): $\mathbf{z}_k, k \in s$.
- Data at an **aggregated** level (number of devices, number of tourists, ...): $\dot{Y}_{s_i}, i \in \mathcal{I}_S$.

- **Work at two steps:**

- From **microdata** \mathbf{z}_k to **partial aggregates** $\dot{Y}_i, i \in \mathcal{I}_S$.
- From **partial aggregates** \dot{Y}_i to **estimators** \hat{Y}_i of target aggregates Y_{U_i} .

Problem

Framework

Concrete proposals

Remarks



The problem: target aggregates

- **Daytime** population size $\rightsquigarrow N_U$
- **Number of** inbound/domestic **tourists** $\rightsquigarrow N_{U_T}$
- **Number of** **commuters** $\rightsquigarrow N_{U_C}$

Problem

Framework

Concrete proposals

Remarks

Methodologically reduced to a problem of **count statistics**

Equivalently reduced to a problem of **population size estimation**



The problem: formulation (I)

Let

- a **finite population** U of interest of **unknown size** N ;
- a partition into **population domains** U_i such that $U = \bigcup_{i \in \mathcal{I}_U} U_i$, mostly **geographical** areas;
- either a **mobile device** (micro) **data** set s or a collection of **partial aggregates** \dot{N}_{S_i} of **count statistics** of interest.

We want to construct:

- an **estimator** \hat{N} for the population size N ;
- an **accuracy measure** for \hat{N} ;
- **estimators** \hat{N}_{U_i} for the population domain sizes N_{U_i} ;
- **accuracy measures** for \hat{N}_{U_i} .

Problem

Framework

Concrete proposals

Remarks



The problem: formulation (II)

Let

- a **finite population** U of interest of **unknown size** N ;
- a partition into **population domains** U_i such that $U = \bigcup_{i \in \mathcal{I}_U} U_i$, mostly **geographical** areas;
- a collection of time intervals t ;
- either a **mobile device** (micro)**data** set s or a time series of collections of **partial aggregates** $\dot{N}_{s_{it}}$ of **count statistics** of interest.

We want to construct:

- a time series of **estimators** \hat{N}_t for the population size N_t in each time interval t ;
- a time series of **accuracy measure** for \hat{N}_t ;
- a time series **estimators** $\hat{N}_{U_{it}}$ for the population domain sizes $N_{U_{it}}$;
- a time series of **accuracy measures** for $\hat{N}_{U_{it}}$.

Problem

Framework

Concrete proposals

Remarks

- Avoid **file drawer problem**
- Assembling of diverse **statistical/AI methods**
 - The **representativeness** problem \rightsquigarrow ecological sampling, Heckman correction, . . .
 - **Geospatial components** in our problem \rightsquigarrow geostatistical modelling, misalignment problem, . . .
 - Role of **official data** as auxiliary information.

Problem

Framework

Concrete proposals

Remarks



Main idea

Our problem is indeed an ecological **animal abundance problem** where the **sampling protocol** is undertaken by the **mobile telecommunication network** and, in some cases (tourists, commuters, . . .), the **processing of microdata** at the mobile device level.

Problem

Framework

Concrete proposals

Remarks



A naive example

Let

- N be the size of the target population U ;
- Y be the number of mobile phone detected in the network;
- p the detection probability of a statistical unit $k \in U$.

We then model $Y \simeq \text{Bin}(N, p)$ so that

- $\hat{N} = \frac{Y}{p}$ is an unbiased estimator for N :

$$\mathbb{E} \left[\frac{Y}{p} \right] = N.$$

- $\hat{\mathbb{V}}(\hat{N}) = Y \cdot \frac{1-p}{p^2}$ is an unbiased estimator for $\mathbb{V}[\hat{N}]$:

$$\mathbb{E} \left[\hat{\mathbb{V}}(\hat{N}) \right] = N \cdot \frac{1-p}{p} = \mathbb{V}(\hat{N}).$$

Problem

Framework

Concrete proposals

Remarks

First elements for the solution

- **Unbiasedness**
- **Variance estimation**
- The population size N to estimate \rightsquigarrow **target process element**
- The detection probability $p \rightsquigarrow$ **observation process element**

Need for more **realistic model** (more elements)

Problem

Framework

Concrete proposals

Remarks



More realistic elements

- p cannot be realistically assumed to be known:

$$p \simeq \text{Model}(\theta).$$

- p can possibly vary from cell to cell in the network:

$$p_i \simeq \text{Model}(\theta_i).$$

- Thus, we are estimating each local partial aggregate

$$(Y_i, p_i) \rightsquigarrow N_{U_i}.$$

Problem

Framework

Concrete proposals

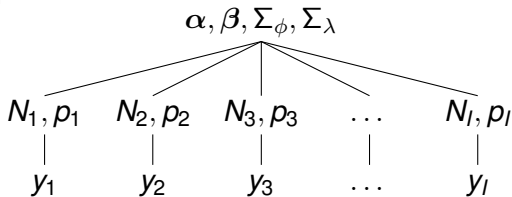
Remarks



Metapopulation
parameters

Local
parameters

Data



A (not so simple) example:

$$\begin{aligned}
 [y_i \mid N_i, p_i] &\simeq \text{Bin}(N_i, p_i) \\
 [N_i \mid \lambda_i] &\simeq \text{Po}(\lambda_i) \\
 [\text{logit}(p_i) \mid \beta, \Sigma_\phi] &\simeq \text{N}(\beta \cdot \mathbf{x}', \Sigma_\phi) \\
 [\lambda_i \mid \alpha, \Sigma_\lambda] &\simeq \text{N}(\alpha, \Sigma_\lambda)
 \end{aligned}$$

General philosophy:

[data|process, parameters][process|parameters][parameters]

Then:

$$n_i = \mathbb{E}[N_i | \text{data}]$$

Problem

Framework

Concrete proposals

Remarks

- **Spatial variations** seem apparently to play a relevant role.
- Spatial variations can be accounted for by using **geostatistical modelling techniques** ($\Sigma_{\phi}, \Sigma_{\lambda}$).

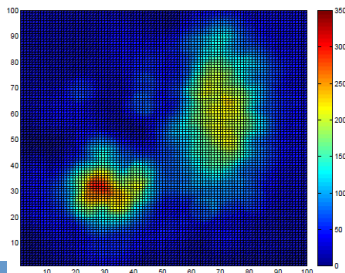
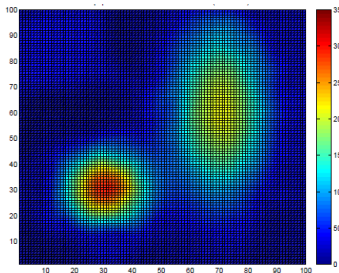
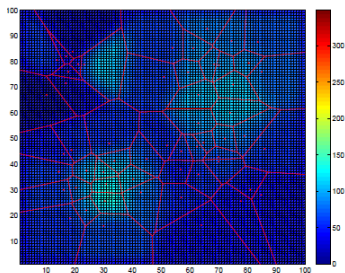
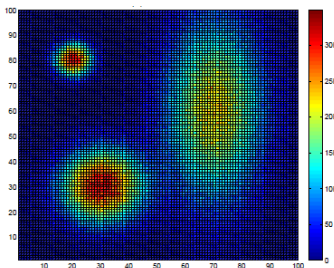
Problem

Framework

Concrete proposals

Remarks





Problem

Framework

Concrete proposals

Remarks

- Calibrate to national official population N^{official}

Example (Deville et al., 2014):

$$\rho_{\text{official}} = \alpha \cdot \sigma_{\text{MD}}^{\beta}, \text{ then } \rho = \frac{\rho^{\text{official}}}{\hat{p}} \cdot \alpha \cdot \sigma_{\text{MD}}^{\beta}.$$

- Do not calibrate to national official population, but only use official figures as covariates \mathbf{x}

Example logit (p_i) = $\beta_0 + \beta_1 \cdot \text{skewness}_i^{\text{income}} + e_i.$

Problem

Framework

Concrete proposals

Remarks



Mobile devices per statistical unit

- Estimates on the **number of statistical units, not of mobile devices**.
- Multiplicity of devices accounted by **penetration rates** τ (see Internal Doc. 3).
- **Model jointly** $\tau_i \cdot p_i$ instead of p_i .

Problem

Framework

Concrete proposals

Remarks



- **Mining microdata sets** in MNOs' hands (with exceptions): restrictions due to data availability (e.g. longitudinal analysis).
- **Internal technical report** by Positium.
- Apparently, mining algorithms are **highly sensitive to domain** (human populations, tourists, mobility–commuters, . . .).

Heckman correction

See presentation by INSSE (B. Oancea).

Problem

Framework

Concrete proposals

Remarks



Admin data methodology

See presentation by INSSE (C. Alexandru).

Problem

Framework

Concrete proposals

Remarks



Inference philosophy

- **Frequentist** vs. **Bayesian** methods in models.
- A priori, Bayesian approach seems **more natural**.
- **Pragmatic** approach \rightsquigarrow **availability of computer tools** at first stage.

Problem

Framework

Concrete proposals

Remarks

