# Smart Statistics for Smart Applications

Book of Short Papers SIS2019

Editors: Giuseppe Arbia, Stefano Peluso, Alessia Pini and Giulia Rivellini

# Preface

# Section 1. Plenary Sessions and Round Table

# Section 2. Invited Papers

# Section 3. Contributed Papers

# Section 4. Posters

# Bayesian estimate of population count with false captures: a latent class approach

## Stima Bayesiana della popolazione con false catture: un approccio basato sulle classi latenti

Davide Di Cecco, Marco Di Zio and Brunero Liseo

**Abstract** We propose a capture–recapture model for estimating the size of a population based on multiple lists in presence of out-of-scope units (false captures). Our Bayesian approach makes use of a class of log–linear models with a latent structure. We also address the presence of sources providing partial information implementing a Gibbs Sampler algorithm which generates a sample from the posterior distribution of the population size in the presence of missing data. The proposed method is applied to simulated data sets.

**Abstract** *Si illustra un metodo di stima per la numerosità di una popolazione basato su liste multiple in presenza di unitá erroneamente conteggiate (false catture). L'approccio bayesiano impiega una sotto classe di modelli log–lineari con struttura latente. Per utilizzare correttamente liste con informazioni parziali, si implementa un algoritmo di tipo Gibbs che genera un campione dalla distribuzione a posteriori del conteggio della popolazione in presenza di dati mancanti. La validità del metodo é verificata attraverso alcune simulazioni.*

**Key words:** Bayesian Analysis, Capture–Recapture, Latent Class

## 1 Introduction

Estimating the size of a population is a central issue in several different fields. Many problems arise in practical applications, because the available information - either

Davide Di Cecco
ISTAT, Rome, e-mail: dicecco@istat.it

Marco Di Zio
ISTAT, Rome, e-mail: dizio@istat.it

Brunero Liseo
Sapienza Università di Roma e-mail: brunero.liseo@uniroma1.it

"capture occasions" in animal abundance problems or "lists" in multiple record systems context - often suffer from *i)* undercoverage, that is, some members of a population have zero chance of being captured or included in some of the lists; *ii)* overcoverage, that is, some units which do not belong to our population are erroneously included ("false captures").

In this note we propose a capture–recapture model which takes into account these two aspects.

There are various proposals in the literature which adopt a complex model to deal with false capture in multiple lists. This methodology is frequent in animal abundance problems: see for example [2], [6], [12]. However, in these works false captures and misidentification are essentially ascribed to record linkage errors. In this work, we want to address a larger class of problems, where false captures do not have an identified source of error. Usually, possible causes of errors are inherently specific to the data at hand and require ad hoc procedures. Obviously, all the available information should be included in the process of identifying the erroneous captures. Ideally, recognizing and deleting spurious cases should constitute a first step of the analysis, followed by a capture–recapture technique on the "cleaned" data. However, the available information is often not sufficient to single out every false capture, and there might remain some uncertainty, difficult to deal with.

To our knowledge, the only contributions dealing with false captures with no restrictive hypothesis on the source of errors in multiple record systems are [8] and [5]. While the former proposes a standard Bayesian log-linear model, the latter extends the approach to include latent variables. However, in both cases, just one list is assumed to suffer from overcoverage.

The problem of overcoverage has received some interest in official statistics, where it is particular relevant when considering administrative sources to estimate population counts. In this context in fact, many possible sources of errors are recognized: record linkage errors are commonly addressed (see e.g., [11]); delays in record updating can lead to erroneous classification of units; more generally the differences in scope of the various sources are hard to harmonize correctly. At present, the most popular approach in official statistics provides for a Dual System Estimator (DSE). All pertinent sources are integrated into a unique population statistical register which is coupled with a survey. An estimate of the overcoverage rate is obtained through the comparison between the (supposedly) error free survey and the register via some supervised model. Then, the estimate is used to "correct" the DSE in some way. The advantage of the procedure is that the DSE is remarkably robust (see, e.g., [1]), and there is no need to rely on any complex model specification.

In our approach, we look at a multiple lists context, in order to exploit the information redundancy. In this context, a series of methodological issues may arise:

- non independence of the captures: it is important to consider potential dependencies among the various lists;
- whereas DSE is robust with respect to violation of basic hypotheses (e.g., the homogeneity of capture probabilities), this is not generally true in Multiple Record Systems;

- some units may have zero probability of being captured in some lists (e.g., lists targeting only specific subpopulation or different periods of time).

The proposal we discuss here relies on the following model assumptions: all possible erroneous captures are defined as random classification errors under a binary model. We assume the presence of two subpopulations: one comprising the out–of–scope units, and the other the in–scope units. Then, a two-component latent class model is assumed to describe the data. In order to model possible dependencies among captures of the same individual in different sources, we relax the classic conditional independence assumption of latent class models and we assume a general log–linear model. In order to deal with the structural absence of some subpopulations from specific lists, we propose to treat the uncatchable units as missing information and develop a missing data approach. This idea has been already proposed in [4]. Here we present its Bayesian counterpart.

## 2 The model

Assume $k$ lists or capture occasions are available, and let $Y_i$ be the random variable indicating whether a unit is included in the $i$–th list, $i = 1, ..., k$ (i.e., has been captured in the $i$–th occasion):

$$Y_i = \begin{cases} 1 & \text{if a unit is captured in the } i\text{–th list;} \\ 0 & \text{otherwise.} \end{cases}$$

Let $\mathbf{Y} = (Y_1, \ldots, Y_k)$ denote the capture profile of a unit, and let $\{P(\mathbf{Y} = \mathbf{y}) = p_{\mathbf{y}}\}_{\mathbf{y} \in \{0,1\}^k}$ be the associated probability distribution.

Let $U(i)$ be the set of units that are catchable by list $i$, and let $U$ be $\bigcup_i U(i)$. Let $U_1$ be our target population, with $U_1 \subset U$. The cardinality of $U$ is $N$, the one of $U_1$ $N_1$. Let $X$ be the latent variable identifying the units belonging to our target population:

$$X = \begin{cases} 1 & \text{if a unit belongs to } U_1; \\ 0 & \text{otherwise.} \end{cases}$$

Let $n_{\mathbf{y}}$ be the number of units having capture profile $\mathbf{y}$, of which $n_{x,\mathbf{y}}$ belong to the latent class $x$ so that $n_{0,\mathbf{y}} + n_{1,\mathbf{y}} = n_{\mathbf{y}}$. The total number of observed unit is $n_{obs}$, while the units having capture history $\mathbf{y} = \mathbf{0} = (0, \ldots, 0)$ are unobserved, so that $\sum_{\mathbf{y} \neq \mathbf{0}} n_{\mathbf{y}} = n_{obs}$, and $N = n_{obs} + n_{\mathbf{0}}$. Note that $n_{1,\mathbf{0}}$ is the number of units in $U_1$ that are not captured, while $n_{0,\mathbf{0}}$ is the number of uncaptured units which are in $U$ but not in $U_1$. We are interested in estimating $N_1 = \sum_{\mathbf{y}} n_{1,\mathbf{y}}$.

The latent class model under the conditional independence assumption (CIA) can be equivalently expressed as the mixture model

$$p_{\mathbf{y}} = \sum_{x=0,1} p_x \prod_{i=1}^{k} p_{y_i|x}, \tag{1}$$

where $p_{y_i|x}$ indicates the conditional probability $P(Y_i = y_i | X = x)$, or as in the log–linear model notation

$$[XY_1][XY_2]\cdots[XY_k], \tag{2}$$

which reports only the higher order interactions (generators) of the model. Any additional interaction term in (2) represents a relaxation of the CIA.

## 2.1 Prior distributions

The usual priors for log–linear models are based on Multivariate Gaussian distributions. Here we propose a different prior based on Dirichlet distributions. We find this approach easier in terms of elicitation of prior knowledge, and also from a computational point of view, since it allows us to develop a Gibbs sampler for obtaining a sample from the posterior distribution of $N_1$, so avoiding the use of a Metropolis–Hastings algorithm.

To illustrate our proposal we start with decomposable models. In this case the prior distribution is simply the product of Dirichlet densities. In [3] it has been demonstrated that, if $G$ is the dependence graph of the decomposable model, $\{\mathscr{C}_1, \ldots, \mathscr{C}_g\}$ are the maximal cliques of $G$, and $(\mathscr{S}_2, \ldots, \mathscr{S}_g)$ are defined as

$$\mathscr{S}_i = \mathscr{C}_i \cap \bigcup_{j=1}^{i-1} \mathscr{C}_j \qquad i = 2, \ldots, g,$$

the joint distribution can be written as the product of conditional distributions:

$$p_G = \prod_{i=1}^{g} p_{\mathscr{C}_i} \left( \prod_{j=2}^{g} p_{\mathscr{S}_j} \right)^{-1} = p_{\mathscr{C}_1} \prod_{i=2}^{g} p_{\mathscr{C}_i | \mathscr{S}_i}, \tag{3}$$

where $p$ over a (sub)graph is the (marginal) distribution over the variables included in the (sub)graph. Let $\Theta$ be the vector of parameters $\Theta = \left( P_{\mathscr{C}_1}, P_{\mathscr{C}_2 | \mathscr{S}_2}, \ldots, P_{\mathscr{C}_g | \mathscr{S}_g} \right)$. We define a prior distribution on $\Theta$ as follows: for each $P_{\mathscr{C}_i | \mathscr{S}_i}$ and for each value of $\mathscr{S}_i$ we set a Dirichlet distribution defined for each possible combination of values $\mathbf{y}_{\mathscr{C}} \in \{0,1\}^{|\mathscr{C}|}$ of the variables in $\mathscr{C}$. The Dirichlet densities are independent by construction, and this class of priors is conjugate to (3).

In the case of a general log–linear model, we made use of the "Bayesian iterative proportional fitting" described in [10] in order to sample from a "Constrained Dirichlet". That is, we generate samples from a Dirichlet distribution which satisfies the constraints given by the log–linear model. This prior has been rarely utilized in literature, and, as far as we know, has never been utilized in capture–recapture analysis.

Regarding $N$, in accordance with the literature on Bayesian capture–recapture, sensible options include:
i) Jeffreys' prior, i.e. $\pi(N) \propto 1/N$;

ii) a hierarchical Poisson prior: $N \sim Poi(\lambda), \quad \lambda \sim Gamma(\alpha, \beta)$;

iii) Rissanen's prior ([9]), $\pi(N) \propto 2^{-\log^*(N)}$, where $\log^*(N)$ is the sum of the positive terms in the sequence $\{\log_2(N), \log_2(\log_2(N)), \ldots\}$.

We further assume that $N$ and $\Theta$ are a priori independent.

## 3 Missing data

We propose a strategy useful to properly include sources which do not operate over certain subpopulations ("incomplete lists"). In fact, if we treat the uncatchable units as sampling zeros, the final population size estimate would be biased.

The idea is to treat the incomplete lists as Missing at Random (MAR) information, i.e. assuming that, if they could operate on the whole population, they would retain the same joint distribution as in the observed subpopulations. In addition, we assume that we can distinguish whether a unit has not been captured in a list by chance or because it is out of the scope of that list, i.e., we can divide the population in strata where different set of lists operates. Then, certain profiles of the captured units are considered as partially observed, and we develop a data augmentation algorithm that imputes the complete capture histories using the rest of the data given the model.

We distinguish completely observed capture profiles, $\mathbf{y}$, from the partially observed capture profiles $\mathbf{y}_{mis}$. In addition, for each stratum, we have a structural zero $\mathbf{z}$ consisting in a different combination of zeros and missing values. For example, in a 4-lists scenario with 2 strata, one where all lists operate and one where the first list does not operate, we have the structural zero $n_{0,0,0,0}$ in the first strata, and $n_{*,0,0,0}$ in the second, where the asterisk denotes the missing information.

Then, our Gibbs algorithm at iteration $t+1$ has the following steps:

1) we sample the components of $\Theta^{(t+1)}$ from their posterior conditional Dirichlet distributions (constrained or not);

2) for each observed $\mathbf{y}$ and $\mathbf{y}_{mis}$, we randomly divide all the observed values $n_{\mathbf{y}}$ and $n_{\mathbf{y}_{mis}}$ into the corresponding consistent complete sequences $n_{x,\mathbf{y}}$ according to their conditional probabilities;

3) if we adopt $\pi(N) \propto 1/N$, it has been demonstrated in [7] that we can sample all structural zero cells counts $n_{\mathbf{z}}$ from a Negative Multinomial distribution. Otherwise, if we choose an informative prior for $N$, we can use a Metropolis-Hasting step to generate a value for $N^{(t+1)}$ and then conditionally sample the structural zero cells such that $\sum_{\mathbf{z}} n_{\mathbf{z}} = N - n_{obs}$;

4) for each generated $n_{\mathbf{z}}$, we sample all complete sequences $n_{x,y}$ consistent with $\mathbf{z}$.

# 4 Simulations

We compare the performance of our procedure in absence and in presence of missing data, and in absence and in presence of prior information. We report the results of a simulation for empirically assessing the proposed algorithm in the various settings. We considered 5 lists, $A$, $B$, $C$, $D$, $E$, and defined two scenarios: in the former, all five sources operate on the whole population, in the latter, there are three strata: one with all sources, another one where 4 out of five sources operate, and one where just three sources operate.

To evaluate the sensitivity to the prior distributions, we considered two additional cases: in the first, we used non informative priors for all parameters (all Dirichlet parameters equal to 1 and $\pi(N) \propto 1/N$), in the second we mimicked an informative context coming from an audit sample: we took a 5% sample of the generated complete population *[XABCDE]*, and fix the parameters of the Dirichlet prior equal to the observed counts in that sample.

We set $N = 10000$ and a proportion of out–of–scope units (both captured and non-captured) equal to 40%, so that the desired total $N_1$ is 6000 in expectation. We generated 500 independent samples from 2 models:

$$\text{Model 1: } [XA][XB][XC][XD][XE], \text{ and Model 2: } [XABC][XD][XE].$$

For each sample, we registered the generated ("true") values of $N_1$ (the target population size), and derived the marginal "observed" counts by omitting the structural zero cells. The model parameters have been set in such a way that the proportion of unobserved units (both in–scope and out–of–scope) is 20% in the scenario without missing information and about 30% in the scenario with three strata. For each simulation we calculated the posterior mean and the 95% credibility interval for $N_1$.

The results are summarized in terms of sample bias in Table 1, and according to the average width of the 95% credibility intervals in Table 2.

**Table 1** Results of the simulations. Sample Bias over the 500 samples in the four scenarios.

|  | Flat prior | | Informative prior | |
|---|---|---|---|---|
|  | No missing | Missings | No missing | Missings |
| Model 1 | 4.1 | -8.3 | 3.6 | 5.1 |
| Model 2 | 3.1 | -12.5 | -4.6 | 6.2 |

Model selection is a critical issue for capture–recapture modeling as population size estimate can be sensitive to changes in the parameterization. To have a hint on the robustness of the procedure under misspecification of the model, in the scenario described above with three strata, we generated a sample from model *[XABC][XD][XE]*, and estimated $N_1$ under two different models: model

**Table 2** Results of the simulations. Average width of the 95% credibility intervals over the 500 samples in the four scenarios.

|  | Flat prior | | Informative prior | |
|---|---|---|---|---|
|  | No missing | Missings | No missing | Missings |
| Model 1 | 309.6 | 531.3 | 258.6 | 378.9 |
| Model 2 | 185.2 | 345.6 | 160.2 | 263.8 |

*[XA][XB][XC][XD][XE]*, and the (non – decomposable) model including all 15 second order interactions but no higher order parameters. Results regarding the second model can be viewed in Figure 1, where one sees that the true value of $N_1$ is comprised in the 95% credibility interval, despite 5 parameters are missing (those relative to *[ABC], [XAB], [XAC], [XBC]* and *[XABC]*).



**Fig. 1** Posterior distributions of $N_1$ under the generating model *[XABC][XD][XE]* (left) and under the all–second–order–interactions model (right). The orange line indicates the true value of $N_1$, the gray area the 95% HPD.

On the converse, the left panel of Figure 2 shows that the estimated posterior distribution of $N_1$ under the CIA model is far from the real value. To evaluate the influence of the prior distributions to compensate for the model misspecification, we set an informative prior as in the previous simulation, where a 5% audit sample establishes the Dirichlet priors parameters. As one can see in the right panel of Figure 2, even though informative priors influence the posterior in the right direction, their contribution seems insufficient to even include the true value of $N_1$ in the credibility interval.

**Fig. 2** Posterior distributions of $N_1$ under the CIA model, flat priors(left) and informative priors (right). The orange line indicates the true value of $N_1$, the gray area the 95% HPD.

# References

1. A. Chao, P.K. Tsay, S.H. Lin, W. Shau, and D. Chao. The applications of capture-recapture models to epidemiological data. *Statistics in medicine*, 20(20):3123–3157, 2001.
2. C. Q. da Silva. Bayesian analysis to correct false-negative errors in capture–recapture photo-ID abundance estimates. *Brazilian Journal of Probability and Statistics*, 23(1):36–48, 2009.
3. A. P. Dawid and S. L. Lauritzen. Hyper markov laws in the statistical analysis of decomposable graphical models. *The Annals of Statistics*, 21(3):1272–1317, 1993.
4. D. Di Cecco, M. Di Zio, D. Filipponi, and I. Rocchetti. Population size estimation using multiple incomplete lists with overcoverage. *J. Off. Stat.*, 34(2):557–572, 2018.
5. D. A. Fegatelli, A. Farcomeni, and L. Tardella. Bayesian population size estimation with censored counts. In *Capture-Recapture Methods for the Social and Medical Sciences*, pages 371–385. Chapman and Hall/CRC, 2017.
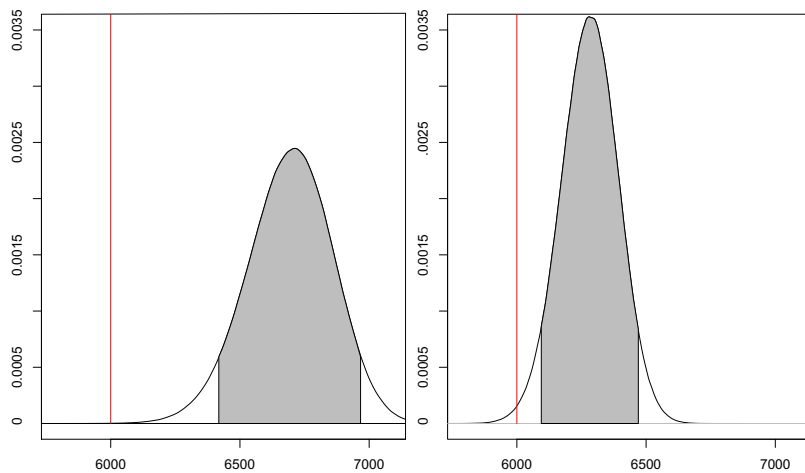6. W. A. Link, J. Yoshizaki, L. L. Bailey, and K. H. Pollock. Uncovering a latent multinomial: analysis of mark–recapture data with misidentification. *Biometrics*, 66(1):178–185, 2010.
7. D. Manrique-Vallier and J. P. Reiter. Bayesian estimation of discrete multivariate latent structure models with structural zeros. *Journal of Computational and Graphical Statistics*, 23(4):1061–1079, 2014.
8. A. M. Overstall, R. King, S. M. Bird, S. J. Hutchinson, and G. Hay. Incomplete contingency tables with censored cells with application to estimating the number of people who inject drugs in Scotland. *Statistics in medicine*, 33(9):1564–1579, 2014.
9. J. Rissanen. A universal prior for integers and estimation by minimum description length. *The Annals of statistics*, 11(2):416–431, 1983.
10. J. L. Schafer. *Analysis of incomplete multivariate data*. CRC press, 1997.
11. A. Tancredi and B. Liseo. A hierarchical Bayesian approach to record linkage and population size problems. *The Annals of Applied Statistics*, 5(2B):1553–1585, 2011.
12. J. A. Wright, R. J. Barker, M. R. Schofield, A. C. Frantz, A. E. Byrom, and D. M. Gleeson. Incorporating genotype uncertainty into mark–recapture-type models for estimating abundance using DNA samples. *Biometrics*, 65(3):833–840, 2009.