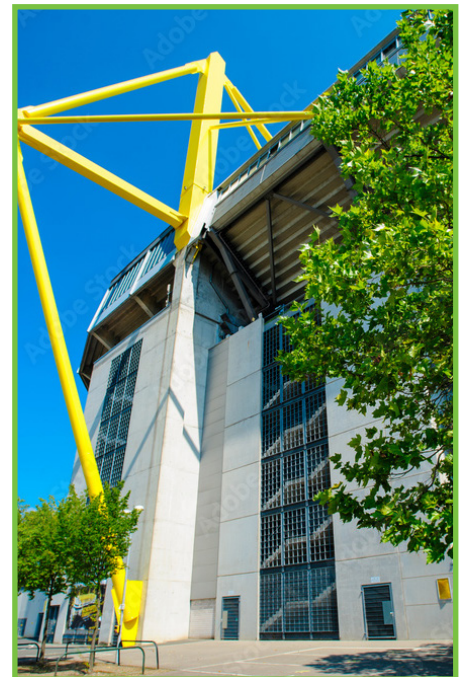technische universität
dortmund

fakultät
statistik

# IWSM 2023

37th International Workshop on Statistical Modelling

16.07. – 21.07.2023

## Dortmund

## Proceedings book

# Proceedings of the 37th International Workshop on Statistical Modelling

**July 17-21, 2023 - Dortmund, Germany**

Editors
Elisabeth Bergherr
Andreas Groll
Andreas Mayr

**Editors:**

**ELISABETH BERGHERR**
University of Göttingen, Chair of Spatial Data Science and Statistical Learning

**ANDREAS GROLL**
TU Dortmund University, Department of Statistics

**ANDREAS MAYR**
University of Bonn, Department of Medical Biometry, Informatics and Epidemiology

## Scientific Committee

Ruggero Bellio
   University of Udine (Italy)
Elisabeth Bergherr (Co-Chair)
   University of Göttingen (Germany)
Fernanda De Bastiani
   University of Pernambuco (Brazil)
María L. Durbán Reguera
   University of Madrid (Spain)
Jan Gertheiss
   Helmut Schmidt University, Hamburg (Germany)
Andreas Groll (Chair)
   TU Dortmund (Germany)
Thomas Kneib
   University of Göttingen (Germany)
Dae-Jin Lee
   IE University, School of Science and Technology, Madrid (Spain)
Andreas Mayr (Co-Chair)
   University of Bonn (Germany)
Fulvia Pennoni
   University of Milano-Bicocca (Italy)
María Xosé Rodríguez Álvarez
   University of Vigo (Spain)
Gunther Schauberger
   TU München (Germany)
Nicola Torelli
   University of Trieste (Italy)
Lola Ugarte
   University of Navarra (Spain)
Nikolaus Umlauf
   University of Innsbruck, (Austria)
Helga Wagner
   University of Linz, (Austria)

# Local Organising Committee

Chiara Balestra
> TU Dortmund University

Elisabeth Bergherr (Co-Host)
> University of Göttingen

Guillermo B. Sánchez
> TU Dortmund University

Jennifer Engel
> TU Dortmund University

Alexander Gerharz
> TU Dortmund University

Colin Griesbach
> University of Göttingen

Andreas Groll (Host)
> TU Dortmund University

Tobias Hepp
> University Erlangen-Nürnberg

Hannah Klinkhammer
> University of Bonn

Andreas Mayr (Co-Host)
> University of Bonn

Hendrik van der Wurp
> TU Dortmund University

# Preface

Dear Participants,

we are more than happy to host the 37th International Workshop on Statistical Modelling in Dortmund, Germany! This is the second year to meet in person after the COVID break, and we hope to have a wonderful time like we did last year in Trieste.

This year we will have 54 contributed talks and more than 60 posters, and it was a tough challenge to pick among the many excellent submissions we had! Thanks again to the scientific committee for putting so much work into the selection process. But we obviously also want to give our thanks to all the researchers, who contributed with their great submissions and made it possible to put together such an excellent set of presentations. Having a special focus on students is a tradition of the Statistical Modelling Society, hence we are especially happy to welcome such a large number of younger researchers contributing to the conference. We are already excited to find out who will win the awards for best student paper, best student presentation and best student poster! The Statistical Modelling Society furthermore awarded travel grants to two students.

We will also have five great invited talks, from different areas in statistics: Brian Reich, Maria Iannario, Alexander Gerharz together with Matthias Kolodziej, Gillian Heller and Simon Wood agreed to give keynotes at the workshop. Furthermore, Andreas Bender and Fabian Scheipl will provide a short course about Piece-wise Exponential (Additive) Models (PEMs / PAMs) before the conference starts.

As always, the IWSM is a one-track conference, leading to a familiar atmosphere and to the possibility for communication between the different fields of statistical modelling.

Looking back at all the years we were participating in great workshops, hosted at so many different universities and all the amazing people we got to meet there, we are both humble and exited to welcome you all to enjoy the conference and your stay at the river Ruhr area with its long tradition of coal mining and steal production, beer brewing and, of course, its omnipresent football vibe.

Andreas Groll, Elisabeth Bergherr and Andreas Mayr

Dortmund, Göttingen and Bonn, July 2023

# Contents

## Part I

## Part II

## Part III

# On the nature of one–inflation in microbial diversity studies

Davide Di Cecco[1], Andrea Tancredi[1]

[1] Sapienza University of Rome, Italy

E-mail for correspondence: `davide.dicecco@uniroma1.it`

**Abstract:** The phenomenon of one–inflation is gaining more and more attention in the recent literature on species abundance and capture–recapture analysis. When analysing frequency count distribution, the excess of singletons is often ascribed to the erroneous inclusion of spurious cases. Various works propose to estimate the true number of singletons relying on the higher, supposedly error–free, counts ("discounting" approach). We argument that, in the case of microbial diversity studies, the generating process of the spurious singletons can be described in terms of false negative record linkage errors. Errors in sequencing the RNA genomes result in chimeric sequences that cannot be associated to the correct species, and constitute missing links that are added to the true singletons. In this scenario, none of the observed frequency counts is assumed to be error–free, and we propose an ABC algorithm to estimate the true frequency counts. The number of true singletons estimated in this way may differ considerably from the discounting approach. This implies different estimates of the diversity as measured, e.g., by Shannon's index. However, curiously, the total population count estimates under the two approaches coincide.

**Keywords:** Species problem; Biodiversity; Linkage Errors; Approximate Bayesian Computing.

## 1 Introduction

The problem of estimating the number of species in a population given a sample arises in many applications in the natural sciences, in linguistics and computer science. Our focus is on applications in microbial ecology. The spread of next generation high-throughput sequencing technology allowed to analyse an unprecedented amount of data on microbial communities. In order to study the biodiversity in a microbial community, an environmental sample is processed to detect, amplify and sequence RNA genomes. The

sequences are clustered into distinct species (or Operational Taxonomic Units) on the basis of a similarity score. The diversity analysis is then conducted on the abundance frequency counts, i.e., the counts $\{n_j\}_{j=1,2,\ldots}$ representing the number of species with $j$ captured occurrences. In most microbial studies, the distribution $\{n_j\}_{j=1,2,\ldots}$ is characterized by an un-expected number of low–abundance species, in particular singletons, accompanied by a low number of very common species. The nature of these singletons has been debated at length, and the presence of spurious single-tons resulting from sequencing errors has been confirmed in various ways (e.g., Quince et al. 2011, Haas et al. 2011). While bioinformatics focuses on avoiding the formation of the so–called chimera sequences, or removing them in a pre–processing step, various statistical contributions attempt to estimate ex–post their number.

The study of one–inflation in frequency count distribution is gaining more and more attention also in the recent capture–recapture literature on hu-man and animal population, which shares many methodological aspects with the species abundance problem, (see, e.g., Godwin and Böhning 2017, Böhning et al. 2019, Tuoto et al. 2022). The possible sources of one–inflation can be categorized as:

- a behavioural effect, where certain units, once captured, avoid subse-quent captures;

- the presence of out-of-scope units, which enter the sample for a pe-culiar error mechanism and should be excluded;

- the presence of missing links in the record linkage procedure employed to create the frequency counts.

Various authors adopted a "discounting" approach to the problem of one–inflation. That is, they propose to ignore the data affected by errors, i.e., the observed singletons, and re-estimate their number on the basis of the counts $n_j$, $j \geq 2$, (see, e.g, Willis and Bunge 2015, Willis 2016, Chiu and Chao 2016). We argument that this approach is consistent for the second mechanism listed above: a model where out-of-scope singletons are added to the baseline distribution of the true counts. We believe that the nature of the spurious cases can alternatively be described by linkage errors. That is, we assume that random errors occurring in sequencing result in the impossibility of a correct classification of the specimen, which cannot be associated to the right existing species. Therefore, we can describe these cases as false negative linkage errors (or missing links), which are added to the true singletons. This approach implies a re–estimation of the "real" frequency counts for all the abundances, not just the singletons. We found that treating the excess of singletons in this way leads to significant differ-ences in the diversity estimates with respect to the discounting approach. In this work we adopt a secondary approach to the linkage problem, i.e., we try to estimate the linkage errors solely on the basis of the vector

$\{n_j\}_{j=1,2,...}$ and our distributional assumptions, as we do not have access to the actual linkage process. Modeling linkage errors in this secondary setting, appears quite complex from a computational point of view. We fix some simplifying assumptions on the type of error in order to tackle the issue, but we still resorted to a Bayesian likelihood–free approach as the most convenient approach.

## 2     One–inflation models

Say we get $n$ species in our sample with abundances $y_1, ..., y_n$, and abundance frequency counts $\{n_j\}_{j\geq1}$. Under an out–of–scope singletons model, the distribution of the abundances (whether the species are observed or not, spurious or not) results in the following mixture of a baseline distribution $\widetilde{f}$ of the non–spurious counts, and a Dirac measure over one:

$$P(Y_i = j \, ; \, \widetilde{f}, \psi) = \begin{cases} (1 - \psi)\widetilde{f}_1 + \psi & \text{if } j = 1; \\ (1 - \psi)\widetilde{f}_j & \text{otherwise,} \end{cases} \qquad (1)$$

where $\psi$ denotes the portion of spurious cases over the total population count. Let $\widetilde{n}_j$ denote the number of species with $j$ non spurious captures. Then, since we assumed $\widetilde{n}_j = n_j$ for $j \geq 2$, we just have to estimate the number of unsampled species $\widetilde{n}_0$, and the number of non–spurious singletons $\widetilde{n}_1$ as a portion of $n_1$. The estimate of the total number of distinct species $\widetilde{N}$ will result as:

$$\sum_{j\geq0} \widetilde{n}_j = \widetilde{n}_0 + n - n_1 + \widetilde{n}_1.$$

A Bayesian estimation of this model presents no difficulties under various parametric families choices for $\widetilde{f}$. A simple Gibbs sampler scheme is the following: under a Beta prior for $\psi$, its posterior is easily updated. Then, a value for $\widetilde{n}_1$ is generated from a Binomial with parameters $1 - \psi$ and $n_1$. Steps to update the values of $\widetilde{n}_0$ and of the parameters of $\widetilde{f}$ are easily found in literature (see, e.g., Tuoto et al. 2022).
Under our missing links proposal, we assume that each sequence has the same probability $\mu$ of being missclassified as a singleton independently from the other. Denote the true number of sampled distinct species as $n^*$, ($n^* < n$). For each species $i$ with $X_i^*$ true captures, we have $M_i$ missing links, such that the registered abundance is reduced from $X_i^*$ to $X_i = X_i^* - M_i$. $M_i$ has the following distribution:

$$P(M_i = m_i \mid X_i^* = x_i^*) = \binom{x_i^*}{m_i}\mu^{m_i}(1 - \mu)^{x_i^* - m_i}, \quad i = 1, ..., n^*. \qquad (2)$$

Let $f^*$ be the baseline distribution of the $X_i^*$. The distribution of the $X_i$ results as a thinning process where a portion $\mu$ of captures disappear. Let

$n_j^*$ denote the true number of species with $j$ captures, and as $N^* = \sum_{j \geq 0} n_j^*$ the total number of distinct species according to the missing links model. Unlike the spurious singletons model, in this case all values $\{n_j^*\}_{j \geq 0}$ have to be estimated, as they will be, in general, different from the observed values. Denote as $\theta$ the parameters defining $f^*$. We adopted an ABC rejection algorithm with the following scheme:

1. generate values for $(\theta, N^*)$ from the priors $\pi(\theta)$ and $\pi(N^*)$;

2. generate values $(n_0^*, n_1^*, n_2^*, ...)$ conditional on $N^*$ and $\theta$;

3. generate a value for $\mu$ from the Beta prior $\pi(\mu)$ (independent from all the rest);

4. generate missing links at random according to the distribution described in (2), given $(n_0^*, n_1^*, n_2^*, ...)$ and $\mu$. Each missing link modifies the observed count, and increments accordingly the number of singletons, thus obtaining the fictitious data $D^*$;

5. retain the current generated values if a measure of distance $\rho$ between the generated data $D^*$ and the observed data $D$ is below a certain threshold $\epsilon$:
$$\rho(D^*, D) < \epsilon.$$

In our application we utilized the euclidean distance.

As the simple ABC rejection scheme can have a low acceptance rate, we further adopted a sequential ABC to accelerate the procedure, as described in Marin et al. 2012.

A simulation study confirmed the correctness of the ABC algorithm under a Poisson, Geometric, and finite mixture of Poisson distributions for $f^*$. Our first finding in a further simulation study comparing the spurious cases and the missing links proposal, has been the substantial identity of the estimates of the total number of species under the two competing models. That is, if we choose $f^*$ and $\widetilde{f}$ in the same family, despite the fact that the estimates of the true abundance frequencies differ under the two models (i.e., $\widetilde{n}_j \neq n_j^*$ for all $j$), we have $N^* = \widetilde{N}$.

To demonstrate this identity, consider the baseline distribution $f^*$ of the values $X_i^*$ introduced above. It is easily demonstrated that, under various parametric family for $f^*$, (notably, if $f^*$ is any mixed Poisson), the distribution of the $X_i$ belongs to the same parametric family. Then, under identifiability of that family, if we use model (1) and take $\widetilde{f}$ in the same family as $f^*$, $\widetilde{f}$ will be identified as the distribution of the $x_i$, for all $x_i > 0$, and $\psi$ would represent the portion of missing links over the total population count. Let $r_0$ be the number of captured species whose occurrences where all missclassified, i.e., such that $M_i = X_i^*$. Let $M$ be the total number of missing links: $M = \sum_{i=1}^{n^*} M_i$. Then we have

$$n^* = n - M + r_0 \quad \text{and} \quad \widetilde{n}_1 = n_1 - M.$$

The missing links mechanism does not affect the number of undetected species $n_0^*$, but under $\widetilde{f}$ the $r_0$ values are included in $\widetilde{n}_0$, i.e., we have $\widetilde{n}_0 = n_0^* + r_0$. Finally, we can write

$$\widetilde{N} = \widetilde{n}_0 + \widetilde{n}_1 + n - n_1 = \widetilde{n}_0 + n - M \ = \ n_0^* + r_0 + n - M = n_0^* + n^* = N^*.$$

As we have said, even if the estimates of the total number of species coincides under the two models, the abundance distribution will differ, and consequently, the estimated diversity will differ. To illustrate the effect of (ignoring) a missing links mechanism on the estimation of diversity, we utilized a simulation study. As a measure of diversity we considered Shannon's diversity $H$ (see, e.g., Chiu and Chao 2016) calculated as:

$$H = \exp\left(-\sum_{j \geq 1} n_j \frac{j}{s} \ln \frac{j}{s}\right). \tag{3}$$

We generated various datasets under Poisson and Geometric baseline distributions, then simulated the effect of missing links to simulate from our model. Then, we estimated Shannon's diversity on the observed data (that is, ignoring any one–inflation mechanism), on the "adjusted" counts as derived from the spurious cases model (that is, trimming the observed number of singletons) and as derived from the ABC procedure for the missing links model. Note that in our Bayesian approach we can easily estimate the posterior distribution of (3). First, we concluded that ignoring an existing one–inflating mechanism, implies a severe overestimation of the diversity. Second, utilizing model (1) when missing links are the true source of error, reduces sensibly the overestimation, but still leads to different results than what can be achieved with an ABC simulating the actual generating process.

## References

Böhning, D., Kaskasamkul, P., van der Heijden, P. G. (2019). A modification of Chao's lower bound estimator in the case of one-inflation. *Metrika*, **82(3)**, 361–384.

Chiu, C. H., Chao, A. (2016). Estimating and comparing microbial diversity in the presence of sequencing errors. *PeerJ*, **4**, e1634.

Godwin, R. T., Böhning, D. (2017). Estimation of the population size by using the one–inflated positive Poisson model. *Journal of the Royal Statistical Society. Series C*, 425–448.

Haas, B.J., Gevers, D., Earl, A.M. et al. (2011) Chimeric 16s rRNA sequence formation and detection in Sanger and 454-pyrosequenced pcr amplicons. *Genome research*, **21(3)**, 494–504.

Marin, J. M., Pudlo, P., Robert, C. P., Ryder, R. J (2012). Approximate Bayesian computational methods. *Statistics and computing*, **22(6)**, 1167-1180.

Quince, C., Lanzen, A., Davenport, R. J. (2011). Removing noise from pyrosequenced amplicons. *BMC bioinformatics*, **12(1)**, $1-18$.

Tuoto, T., Di Cecco, D., Tancredi, A. (2022). Bayesian analysis of one-inflated models for elusive population size estimation. *Biometrical Journal*, **64(5)**, $912-933$.

Willis, A., Bunge, J. (2015). Estimating diversity via frequency ratios. *Biometrics*, **71(4)**, $1042-1049$.

Willis, A. (2016). Species richness estimation with high diversity but spurious singletons. *arXiv preprint arXiv:1604.02598*.