

Contaminated Model for Species Sampling

M. Barin V. Giorgi A. Grignani
C. Laurenti Argento F. E. Quadri M. Tracanella

Politecnico di Milano

Tutor: Riccardo Corradin

November 11, 2022

Species Sampling

1 IDEA

- A random sample taken from a population.

2 WHAT ARE THE PROBLEMS TO SOLVE

- how do we estimate the number of species in the population?
- how do we estimate the probability of discovering a new species in one additional sample?

Remarkable Examples

- Microbiological studies
- Linguistic
- Genetics
- Observational ecological studies

It is common that data are contaminated

An additional problem: Contamination

Contaminated Model

- ① Contamination derives from human errors: frequency one.
- ② Rare Events can also show up with frequency one.
- ③ Our goal is distinguish between contamination and events.

Accounting for the contamination in the model requires to generalize the models commonly used in literature

Gibbs-type Prior

Let $\{X_i\}_{i \geq 1}$ be a sequence of exchangeable observations. We call this sequence a species sampling sequence if there exists a random probability measure \tilde{p} such that $X_i \sim \tilde{p}$, where:

$$\tilde{p} = \sum_{j \geq 1} p_j \delta_{Z_j} + \left(1 - \sum_{j \geq 1} p_j\right) P_0$$

- $\{p_j\}_{j \geq 1}$: sequence of random weights
- $\{Z_j\}_{j \geq 1}$: sequence of random atoms
- P_0 : contaminant (diffuse) probability measure

We define \tilde{p} a non-proper prior when: $\sum_{j \geq 1} p_j < 1$

The Non-proper Prior

- ① We will show how non-proper models are particularly suited to take into account contaminated observations or more generally observations with frequency one.
- ② Remember: the component P_0 generates singleton blocks, which are called dust in the probabilistic literature on random partitions.

Benefits

- ① Takes into account observations with frequency 1.
- ② Their predictive structure (w.r.t. the non contaminated prior).
- ③ Maintains the analytical tractability of non contaminated priors.

Contaminated Gibbs-type Priors

- 1 Gibbs-type priors are predominant priors in species sampling problems.
- 2 We introduce a new subfamily by including a contaminant component.
- 3 We will call this family: Contaminated Gibbs-type priors.

Definition

$$\tilde{p} = \beta \tilde{q} + (1 - \beta) P_0 \quad \beta \in [0, 1]$$

- \tilde{p} : contaminated Gibbs-type prior
- \tilde{q} : Gibbs-type prior
- P_0 : contaminant probability measure
- β : weight

Our Path

- ① We hope to show the advantages of using these types of priors compared to the non contaminated ones.
- ② One of the main advantages is that: it will become less probable to make wrong conclusions.
- ③ To validate our analytical conclusion regarding contaminated Gibbs priors, we will use the Global Biodiversity Information Facility dataset.