

Modeling incomplete compositional datasets

Andriette Bekker¹, Jason Pillay¹, Antonio Punzo², and
Cristina Tortora³

¹ *University of Pretoria, South Africa*

² *University of Catania, Italy*

³ *San Jose State University, California, USA*

Abstract

Incomplete compositional data analysis is hindered by the lack of likelihood-based methods that directly handle missing proportions on the simplex. This paper introduces an estimation technique that operates directly on the original data, thereby accommodating missing values while preserving the interpretability of the variables within their natural domain. Simulation studies demonstrate robust performance in parameter estimation and imputation across missingness mechanisms. To showcase its practicality, the proposed algorithm is applied to two real datasets exhibiting distinct missingness patterns. Analysis of the American Time Use dataset (with values missing at random) identifies six distinct categories of daily activities, yielding interpretable profiles of daily activity patterns and the demographics associated with them. Application to PM_{2.5} species data from the Air Quality System, supports a four-component mixture model, providing a descriptive profile of measured PM_{2.5} speciation across the United States.

Keywords

Compositional data, Missing data imputation, Mixture models
Mixture models