

# Modeling of compositional data: a multilevel approach to benthic cover Abrolhos bank.

Pamela M. Chiroque-Solano

Multidisciplinary Institute, Federal Rural University of Rio de Janeiro and  
Institute of Biology and SAGE-COPPE, Federal University of Rio de Janeiro, Brazil.

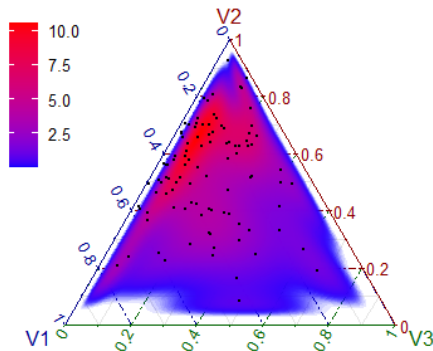
13th of October, 2021

# Framework

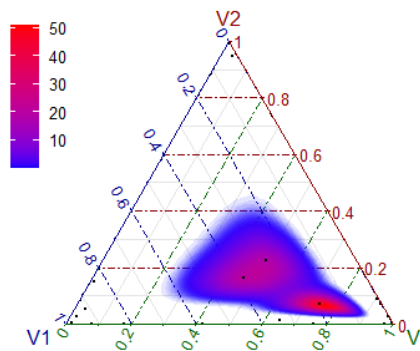
- Multivariate regression with constrained response.
- Challenge:
  - ▶ Unbalanced;
  - ▶ Lot of missing data;
  - ▶ Identificability issues

Objective: To study the variability by site

Case 1:  $\log \phi = 1$



Case 3:  $\log \phi = -1$



**Figure:** Consider a three-components (compositional data). The first simplex contains the information for high entropy (case 1), and low entropy (case 3).

# Model

## Maier (2014) and Holger (2018)

Filtered information through the decomposition of  $\alpha$

- $\mathbf{Y}_l \sim D(\mu_l, \phi_l)$  with parameter  $\alpha_{cl} = \mu_{cl}\phi_l$
- $\mu_{cl}$  : level term
- $\phi_l$  : precision term

## Reference component: $c^*$

- Alternative parametrization:  $c^*$  should be chosen.
- Stochastic representation for Dirichlet random vector

## Sharing information equation

$$\begin{aligned}\beta_{cl} &= \beta_c + \epsilon_{\beta_l}, & \epsilon_{\beta_l} &\sim \mathcal{N}(0, V_\beta) \\ \theta_l &= \theta + \epsilon_{\theta_l}, & \epsilon_{\theta_l} &\sim \mathcal{N}(0, V_\theta)\end{aligned}$$

# Inference procedure

Let  $\Theta = (\beta, \phi)$  be the vector of parameters

Proper independent prior distribution for the parametric vector  $\Theta$  are Normal with zero mean and precision  $1/K$  for all effects of the model.

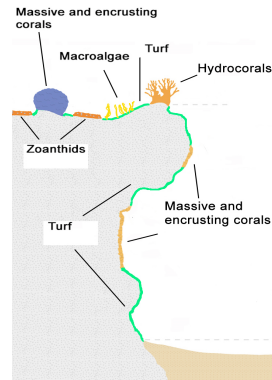
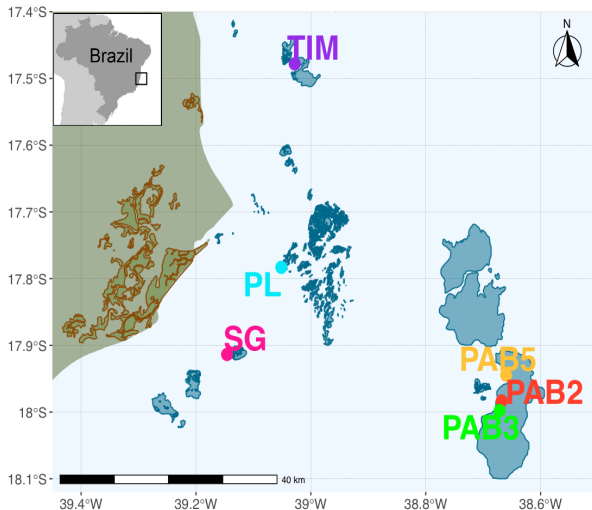
The joint posterior distribution does not have a known closed form

$$\pi(\Theta \mid \mathbf{y}) \propto L(\Theta \mid \mathbf{y}) \prod_l^L \pi(\phi_l) \prod_c^C \pi(\beta_{cl}) \quad (1)$$

Sampling from the posterior distribution

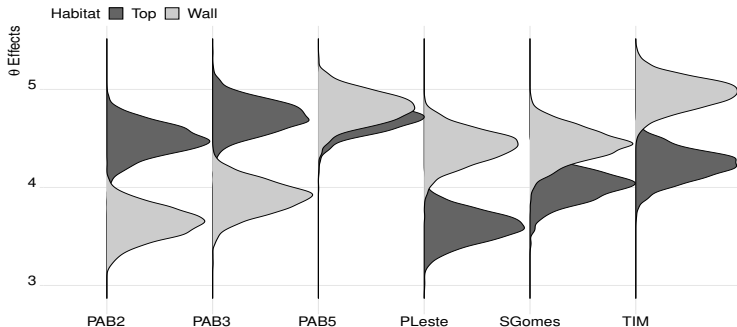
by Markov chain Monte Carlo (MCMC) via the Stan software.

# Application



# Results

**Figure:** Posterior density of the  $\theta$  effect for each of the nine components by site and habitat levels.



The results validate the original hypotheses

Sites near the coast (inshore) are more variable than the offshore sites.

# Conclusions and Future Work

## Main conclusions

- The proposed model quantifies the heteroscedasticity through precision effects via hierarchical structures by site;
- The method is flexible;
- The reference component has been chosen using objective criteria;
- The proposal allows to obtain adequate predictions.
- This work contributes to the United Nations's Sustainable Development Goal 14 - "Life Under Water".



# References

- Gelman, Andrew, and Jennifer Hill. 2006. Data Analysis Using Regression and Multilevel/Hierarchical Models. Analytical Methods for Social Research. Cambridge University Press.
- Holger, and Sennhenn-Reulen. 2018. “Bayesian Regression for a Dirichlet Distributed Response Using Stan.”
- Maier, Marco J. 2014. “DirichletReg: Dirichlet Regression for Compositional Data in R.” Research Report Series/Department of Statistics and Mathematics 125. Vienna: WU Vienna University of Economics and Business.
- Wang, K., G. Tian, and M Tang. 2011. Dirichlet and Related Distributions: Theory, Methods and Applications. Wiley Series in Probability; Statistics.

Thank you

- Guido A. Moreira, UMinho.
- Marine Biodiversity and Conservation Lab at UFRJ.
- The Fundação Espírito Santense de Tecnologia, FEST.
- The Fundação de Amparo à Pesquisa do Estado do Rio de Janeiro, FAPERJ - E-26/200.016/2021.

Contact

[pamela@ufrj.br](mailto:pamela@ufrj.br)