

Template to prepare preprints and manuscripts using markdown and github actions

Timothée Poisot ^{1,2}

¹ Université de Montréal ² Québec Centre for Biodiversity Sciences

Correspondance to:

Timothée Poisot — timothee.poisot@umontreal.ca

This work is released by its authors under a CC-BY 4.0 license



Last revision: *December 7, 2021*

prediction of species interaction networks is gaining momentum as a way to circumvent limitations in data volume. Yet, ecological networks are challenging to predict because they are typically small and sparse. Dealing with extreme class imbalance is a challenge for most binary classifiers, and there are currently no guidelines as to how predictive models can be trained. simple mathematical arguments and numerical experiments in which a variety of classifiers (for supervised learning) are trained on simulated networks, we develop a series of guidelines related to the choice of measures to use for model selection, and the degree of unbiasing to apply to the training dataset. accuracy and the ROC-AUC are not informative measures for the performance of interaction prediction. PR-AUC is a fairer assesment of performance. In some cases, even standard measures can lead to selecting a more biased classifier because the effect of connectance is strong. The amount of correction to apply to the training dataset depends as a function of the classifier and the network connectance. results reveal that training machines to predict networks is a challenging task, and that in virtually all cases, the composition of the training set needs to be experimented on before performing the actual training. We discuss these consequences in the context of the low volume of data.

- 1 example on diagnostic test: rare events are hard to detect even with really good models
- 2 summary of model challenges for networks
- 3 list of problems to solve - baseline values and response to bias - effect of training set bias on performance -
- 4 which models need the least amount of interactions to work
- 5 summary of the results

6 **Baseline values**

7 **Numerical experiments**

8 **Effect of training set on performance**

9 **Required amount of positives to get the best performance**

10 **Guidelines for prediction**

11 **References**