

# Guidelines for the supervised learning of species interactions

Timothée Poisot<sup>1,2</sup>

<sup>1</sup> Université de Montréal; <sup>2</sup> Québec Centre for Biodiversity Sciences

## Correspondance to:

Timothée Poisot — timothee.poisot@umontreal.ca

1. The 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.
2. Using 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.
3. Classifier 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.
4. These 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.

## Keywords:

species interaction networks  
binary classifiers  
machine learning  
regression  
supervised learning

example on diagnostic test: rare events are hard to detect even with really good models

summary of model challenges for networks - Strydom et al. (2021) importance of drawing on traits + validation is challenging - Whalen et al. (2021) machine learning from genomics

Binary classifiers are usually assessed by measuring properties of their confusion matrix, *i.e.* the contingency table reporting true/false positive/negative hits. A confusion matrix is laid out as

$$\begin{pmatrix} tp & fp \\ fn & tn \end{pmatrix},$$

wherein  $tp$  is the number of interactions predicted as positive,  $tn$  is the number of non-interactions predicted as negative,  $fp$  is the number of non-interactions predicted as positive, and  $fn$  is the number of interactions predicted as negative. Almost all measures based on the confusion matrix express rates of error or success as proportions, and therefore the values of these components matter in a *relative* way. At a coarse scale, a classifier is *accurate* when the trace of the matrix divided by the sum of the matrix is close to 1, with other measures focusing on different ways in which the classifier is wrong.

The same approach is used to evaluate *e.g.* species distribution models (SDMs). Indeed, the training and evaluation of SDMs as binary classifiers suffers from the same issue of low prevalence. In a previous work, Allouche et al. (2006) suggested that  $\kappa$  was a better test of model performance than the True Skill Statistic (TSS), which we will refer to as Youden's informedness (or  $J$ ); these conclusions were later criticized by Somodi et al. (2017), who emphasized that informedness' relationship to prevalence depends on assumptions about bias in the model. Although this work offers recommendations about the comparison of models, it doesn't establish baselines or good practices for training on imbalanced ecological data. Within the context of networks, there are three specific issues that need to be addressed. First, what values of performance measures are we expecting for a classifier that has poor performance? This is particularly important as it can evaluate whether low prevalence can lull us into a false sense of predictive accuracy. Second, independently of the question of model evaluation, is low prevalence an issue for *training*, and can we remedy it? Finally, because the low amount of data on interaction makes a lot of imbalance correction methods (see *e.g.* Branco et al., 2015) hard to apply, which indicators can be optimized with the least amount of positive interaction data?

We establish that due to the low prevalence of interactions, even poor classifiers applied to food web data will reach a high accuracy; this is because the measure is dominated by the accidental correct predictions of negatives. The  $F_1$  score and positive predictive values are less sensitive to bias, but **TODO**

Chicco et al. (2021) - MCC maximizes other measures, other measures do not maximize MCC, except notably when prevalence is low, or the baseline guessing level is uncertain, which applies to networks. In this cases, informedness should be used as a maximization criterion. Formulating guidelines is particularly important, because most of the literature existing on optimizing classifier performance in the life sciences is focused on genomics applications (which has very specific challenges, see a recent discussion by Whalen et al., 2021), and can give contradictory recommendations (Boughorbel et al., 2017; Chicco & Jurman, 2020; Delgado & Tibau, 2019). This points not to a deficiency in the literature, but rather to a need for domain-specific evaluation of how the particular ways in which datasets are biased can affect the performance of predictive models.

1

## Baseline values

**1.1. Confusion matrix with skill and bias** In this section, we will assume a network of connectance  $\rho$ , *i.e.* having  $\rho S^2$  interactions (where  $S$  is the species richness), and  $(1 - \rho)S^2$  non-interactions. Therefore, the vector describing the *true* state of the network is a column vector  $\mathbf{o}^T = [\rho(1 - \rho)]$  (we can safely drop the  $S^2$  terms, as we will work on the confusion matrix, which ends up expressing *relative* values).

In order to write the values of the confusion matrix for a hypothetical classifier, we need to define two characteristics: its skill, and its bias. Skill, here, refers to the propensity of the classifier to get the correct answer (*i.e.* to assign interactions where they are, and to not assign them where they are not). A no-skill classifier guesses at random, *i.e.* it will guess interactions with a probability  $\rho$ . The predictions of a no-skill classifier can be expressed as a row vector  $\mathbf{p} = [\rho(1 - \rho)]$ . The confusion matrix  $\mathbf{M}$  for a no-skill classifier is given by the element-wise product of these vectors  $\mathbf{o} \odot \mathbf{p}$ , *i.e.*

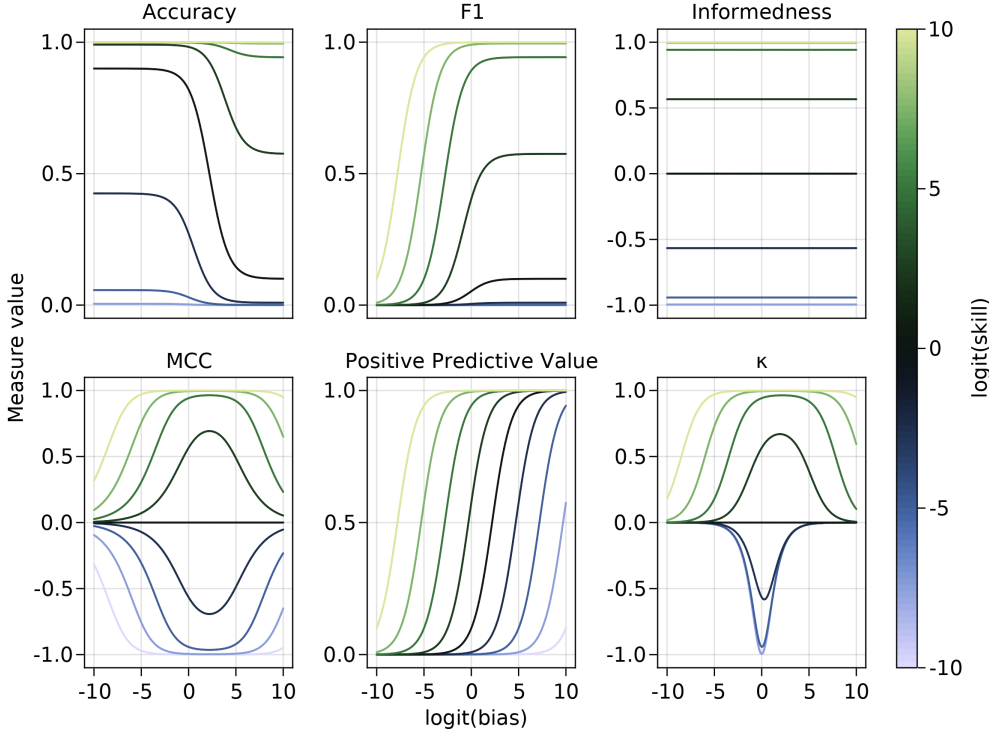
$$\mathbf{M} = \begin{pmatrix} \rho^2 & \rho(1 - \rho) \\ (1 - \rho)\rho & (1 - \rho)^2 \end{pmatrix}.$$

In order to regulate the skill of this classifier, we can define a skill matrix  $\mathbf{S}$  with diagonal elements equal to  $s$ , and off-diagonal elements equal to  $(1 - s)$ , and re-express the skill-adjusted confusion matrix as  $\mathbf{M} \odot \mathbf{S}$ , *i.e.*

$$\begin{pmatrix} \rho^2 & \rho(1 - \rho) \\ (1 - \rho)\rho & (1 - \rho)^2 \end{pmatrix} \odot \begin{pmatrix} s & (1 - s) \\ (1 - s) & s \end{pmatrix}.$$

Note that when  $s = 0$ ,  $\text{Tr}(\mathbf{M}) = 0$  (the classifier is *always* wrong), when  $s = 0.5$ , the classifier is no-skill and guesses at random, and when  $s = 1$ , the classifier is perfect.

The second element we can adjust in this hypothetical classifier is its bias, specifically its tendency to over-predict interactions. Like above, we can do so by defining a bias matrix  $\mathbf{B}$ , where interactions are over-predicted with probability  $b$ , and express the final classifier confusion matrix as  $\mathbf{M} \odot \mathbf{S} \odot \mathbf{B}$ , *i.e.*



**Figure 1** Consequences of changing the classifier skills ( $s$ ) and bias ( $b$ ) for a connectance  $\rho = 0.15$ , on accuracy,  $F_1$ , positive predictive value, and  $\kappa$ . Accuracy increases with skill, but also increases when the bias tends towards estimating *fewer* interactions. The  $F_1$  score increases with skill but also increases when the bias tends towards estimating *more* interactions; PPV behaves in the same way. Interestingly,  $\kappa$  responds as expected to skill (being negative whenever  $s < 0.5$ ), and peaks for values of  $b \approx 0.5$ ; nevertheless, the value of bias for which  $\kappa$  is maximized is *not*  $b = 0.5$ , but instead increases with classifier skill. In other words, at equal skill, maximizing  $\kappa$  would lead to select a *more* biased classifier.

$$\begin{pmatrix} \rho^2 & \rho(1-\rho) \\ (1-\rho)\rho & (1-\rho)^2 \end{pmatrix} \odot \begin{pmatrix} s & (1-s) \\ (1-s) & s \end{pmatrix} \odot \begin{pmatrix} b & b \\ (1-b) & (1-b) \end{pmatrix}.$$

The final expression for the confusion matrix in which we can regulate the skill and the bias is

$$\mathbf{C} = \begin{pmatrix} s \times b \times \rho^2 & (1-s) \times b \times \rho(1-\rho) \\ (1-s) \times (1-b) \times (1-\rho)\rho & s \times (1-b) \times (1-\rho)^2 \end{pmatrix}.$$

In all further simulations, the confusion matrix  $\mathbf{C}$  is transformed so that it sums to 1.

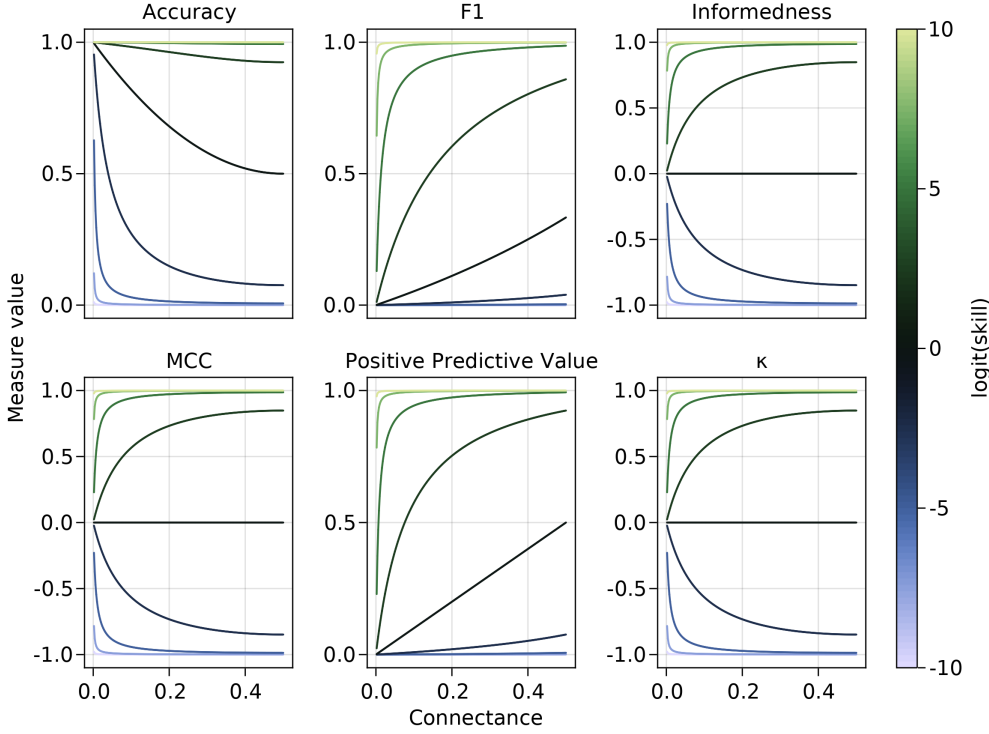
**1.2. What are the baseline values of performance measures?** In this section, we will change the values of  $b$  and  $s$  for a given value of  $\rho$ , and see how the values of common performance measures for binary classification are affected. Specifically, we will focus on four quantities: the accuracy  $((\text{tp} + \text{tn})/(\text{tp} + \text{tn} + \text{fp} + \text{fn}))$ , the balanced accuracy  $(\text{tp}/(2(\text{tp} + \text{fn})) + \text{tn}/(2(\text{tn} + \text{fp})))$ , Youden's  $J$   $(\text{tp}/(\text{tp} + \text{fn}) + \text{tn}/(\text{tn} + \text{fp}) - 1)$ , and the  $F_1$  score  $(2\text{tp}/(2\text{tp} + \text{fp} + \text{fn}))$ .

**Justification** of why these 4

Assuming a no-skill unbiased classifier ( $s = 0.5$ , i.e.  $\mathbf{C} = \mathbf{M}$  after normalization), the accuracy is  $\rho^2 + (1-\rho)^2$ , the balanced accuracy is 0.5, Youden's  $J$  is 0, and  $F_1 = \rho$ . In other words, given a connectance  $\rho = 0.05$ , we expect that a classifier guessing at random would still achieve an accuracy of 0.905; for a connectance of  $\rho = 0.01$ , this accuracy *increases* to over 0.98. In other words, networks with fewer interactions have inherently higher accuracy, because it is easy to predict the overwhelming majority of non-interactions right.

In order to examine how these values change w.r.t. the skill and bias, we performed a grid exploration of the values of  $\text{logit}(s)$  and  $\text{logit}(b)$  linearly from  $-10$  to  $10$ , and visualize the result for a connectance of  $\rho = 0.15$ , which is within the range of usually observed connectance values for empirical food webs.

**1.3. Are the measures affected by connectance?**



**Figure 2** TODO

2

## Numerical experiments

In the following section, we will generate random networks, and train four binary classifiers (as well as an ensemble model using the sum of the outputs) on 30% of the interaction data. Networks are generated by picking random generality  $g$  and vulnerability  $v$  traits for  $S = 200$  species uniformly on the unit interval, and assigning an interaction from species  $i$  to species  $j$  if  $0.2g_i - \xi \leq v_j \leq 0.2g_i + \xi$ , where  $\xi$  is a constant regulating the connectance of the networks, and varies uniformly in  $[5 \times 10^{-3}, 10^{-1}]$ . This model gives fully interval networks that are close analogues to the niche model (Williams & Martinez, 2000), but has the benefit of only relying on two features ( $g_i, v_j$ ), and having the exact same rule for all interactions. It is, therefore, a simple case which most classifiers should be able to learn.

The training sample is composed of 30% of the  $4 \times 10^4$  possible entries in the network, *i.e.*  $n = 12000$ . Out of these interactions, we pick a proportion  $\nu$  (the training set bias) to be positive, so that the training set has  $\nu n$  interactions, and  $(1 - \nu)n$  non-interactions. We vary  $\nu$  uniformly in  $]0, 1[$ . This allows to evaluate how the measures of binary classification performance respond to artificially rebalanced dataset for a given network connectance. Note that both  $\xi$  and  $\nu$  are sampled from a distribution rather than being picked on a grid; this is because there is no direct relationship between the value of  $\xi$  and the connectance of the simulated network, and therefore the precise value of  $\xi$  is not relevant for the analysis of the results.

The dataset used for numerical experiments is composed of 20000 such  $(\xi, \nu)$  pairs, on which four learners are trained: a decision tree regressor, a boosted regression tree, a ridge regressor, and a random forest regressor. All models were taken from the `MLJ.jl` package (Blaom et al., 2020; Blaom & Vollmer, 2020) in Julia 1.7 (Bezanson et al., 2017). In order to pick the best adjacency matrix for a given learner, we performed a thresholding approach using 500 steps on predictions from the testing set, and picking the threshold that maximized Youden's informedness, which is usually the optimized target for imbalanced classification. During the thresholding step, we measured the area under the receiving-operator characteristic (ROC-AUC) and precision-recall (PR-AUC) curves, as measures of overall performance over the range of returned values. We report the ROC-AUC and PR-AUC, as well as a suite of other measures as introduced in the next section, for the best threshold. The ensemble model was generated by summing the predictions of all component models on the testing set (ranged in  $[0, 1]$ ), then put through the same

thresholding process. The complete code to run the simulations is given as an appendix.

After the simulations were completed, we removed all runs (*i.e.* pairs of  $\xi$  and  $\nu$ ) for which at least one of the following conditions was met: the accuracy was 0, the true positive or true negative rates were 0, the connectance was larger than 0.2. This removes both the obviously failed model runs, and the networks that are more densely connected compared to the connectance of empirical food webs (and are therefore less difficult to predict, being less imbalanced).

## 2.1. Effect of training set bias on performance

## 2.2. Required amount of positives to get the best performance

3

---

## Guidelines for prediction

---

## References

- Allouche, O., Tsoar, A., & Kadmon, R. (2006). Assessing the accuracy of species distribution models: Prevalence, kappa and the true skill statistic (TSS). *Journal of Applied Ecology*, 43(6), 1223–1232. <https://doi.org/10.1111/j.1365-2664.2006.01214.x>
- Bezanson, J., Edelman, A., Karpinski, S., & Shah, V. (2017). Julia: A Fresh Approach to Numerical Computing. *SIAM Review*, 59(1), 65–98. <https://doi.org/10.1137/141000671>
- Blaom, A. D., Kiraly, F., Lienart, T., Simillides, Y., Arenas, D., & Vollmer, S. J. (2020). MLJ: A Julia package for composable machine learning. *Journal of Open Source Software*, 5(55), 2704. <https://doi.org/10.21105/joss.02704>
- Blaom, A. D., & Vollmer, S. J. (2020, December 31). *Flexible model composition in machine learning and its implementation in MLJ*. <http://arxiv.org/abs/2012.15505>
- Boughorbel, S., Jarray, F., & El-Anbari, M. (2017). Optimal classifier for imbalanced data using Matthews Correlation Coefficient metric. *PloS One*, 12(6), e0177678. <https://doi.org/10.1371/journal.pone.0177678>
- Branco, P., Torgo, L., & Ribeiro, R. (2015, May 13). *A Survey of Predictive Modelling under Imbalanced Distributions*. <http://arxiv.org/abs/1505.01658>
- Chicco, D., & Jurman, G. (2020). The advantages of the Matthews correlation coefficient (MCC) over F1 score and accuracy in binary classification evaluation. *BMC Genomics*, 21(1), 6. <https://doi.org/10.1186/s12864-019-6413-7>
- Chicco, D., Tötsch, N., & Jurman, G. (2021). The Matthews correlation coefficient (MCC) is more reliable than balanced accuracy, bookmaker informedness, and markedness in two-class confusion matrix evaluation. *BioData Mining*, 14, 13. <https://doi.org/10.1186/s13040-021-00244-z>
- Delgado, R., & Tibau, X.-A. (2019). Why Cohen's Kappa should be avoided as performance measure in classification. *PloS One*, 14(9), e0222916. <https://doi.org/10.1371/journal.pone.0222916>
- Somodi, I., Lepesi, N., & Botta-Dukát, Z. (2017). Prevalence dependence in model goodness measures with special emphasis on true skill statistics. *Ecology and Evolution*, 7(3), 863–872. <https://doi.org/10.1002/ece3.2654>
- Strydom, T., Catchen, M. D., Banville, F., Caron, D., Dansereau, G., Desjardins-Proulx, P., Forero-Muñoz, N. R., Higino, G., Mercier, B., Gonzalez, A., Gravel, D., Pollock, L., & Poisot, T. (2021). A roadmap towards predicting species interaction networks (across space and time). *Philosophical Transactions of the Royal Society B: Biological Sciences*, 376(1837), 20210063. <https://doi.org/10.1098/rstb.2021.0063>
- Whalen, S., Schreiber, J., Noble, W. S., & Pollard, K. S. (2021). Navigating the pitfalls of applying machine learning in genomics. *Nature Reviews Genetics*, 1–13. <https://doi.org/10.1038/s41576-021-00434-9>

Williams, R., & Martinez, N. (2000). Simple rules yield complex food webs. *Nature*, 404, 180–183.  
<http://userwww.sfsu.edu/>