

1 **Appendix S3: Derivation of D: ‘Deviation from the best model’**

2 Model predictions from Simulation 1 and Simulation 2 were evaluated using a novel metric derived
3 from first principles that we named ‘deviation from the best model’ (D). This metric compares the
4 probability of obtaining the true (raw) abundance (i.e. before sampling occurs) in each cell based on the
5 model prediction, with the probability of obtaining a prediction equal to the true (raw) abundance i.e.
6 predictions produced by a perfect model (Eq. 1). For each grid cell i , the probability of obtaining the
7 true (raw) abundance (pA_i) was estimated from a Poisson probability distribution with a mean equal to
8 the predicted mean abundance (\bar{A}_i) for that cell. The summed natural logs of these probabilities across
9 the study area represents the overall probability of obtaining the true (raw) abundances under the model
10 predictions. This is then expressed as a ratio against the summed natural log probabilities for each cell
11 ($q\bar{A}_i$) that would be obtained for a perfect model where the true (raw) abundance is equal to the predicted
12 mean abundance.

13
14 Eq. 1
$$D = \frac{\sum(\ln(pA_i|\bar{A}_i))}{\sum(\ln(q\bar{A}_i|\bar{A}_i))}$$

15

16 The rationale behind creating a new evaluation metric is that our generation of occurrence points was
17 based on a Poisson process and was deliberately zero-inflated, so there are an extremely large proportion
18 of 0’s and 1’s, and the highest ‘abundance count’ is only six. Therefore, a metric based on the probability
19 of obtaining the raw data from the model, rather than a direct assessment of the actual values, would
20 provide a more appropriate measure of model performance and fit that was not as weighted by the large
21 proportion of zeroes in the data. Using traditional binary presence-absence classification metrics would
22 results in penalties against large predictions of abundance in comparison to the true raw abundance: the
23 observed values of 0 and 1 which are most common in our dataset would score more highly using binary
24 classification methods. The magnitude of the difference between prediction and true (raw) abundance
25 observations will scale with the mean, leading to greater “inaccuracy” in cells with large numbers of
26 individuals, regardless of the model used. Therefore, a metric that considers the actual abundance value
27 rather than just presence-absence is likely to provide a more accurate assessment of model predictive
28 power.