

---

# APPLYING THE MAXIMUM ENTROPY PRINCIPLE TO NEURAL NETWORKS ENHANCES MULTI-SPECIES DISTRIBUTION MODELS

---

A PREPRINT

Maxime Ryckewaert<sup>1</sup>, Diego Marcos<sup>1</sup>, Christophe Botella<sup>1</sup>, Maximilien Servajean<sup>2</sup>, Pierre Bonnet<sup>3</sup> and Alexis Joly<sup>1</sup>

<sup>1</sup>Inria, Univ Montpellier, Montpellier, France

<sup>2</sup>LIRMM, AMIS, Univ Paul Valéry Montpellier, Univ Montpellier, CNRS, Montpellier, France

<sup>3</sup>AMAP, Univ Montpellier, CIRAD, CNRS, INRAE, IRD, Montpellier, France

January 16, 2025

## ABSTRACT

The rapid expansion of citizen science initiatives has led to a significant growth of biodiversity databases, and particularly presence-only (PO) observations. PO data are invaluable for understanding species distributions and their dynamics, but their use in a Species Distribution Model (SDM) is curtailed by sampling biases and the lack of information on absences. Poisson point processes are widely used for SDMs, with Maxent being one of the most popular methods. Maxent maximises the entropy of a probability distribution across sites as a function of predefined transformations of variables, called features. In contrast, neural networks and deep learning have emerged as a promising technique for automatic feature extraction from complex input variables. Arbitrarily complex transformations of input variables can be learned from the data efficiently through backpropagation and stochastic gradient descent (SGD). Yet, deep learning was mainly developed for classification problems, and learning robust features and species abundances across space while properly correcting for sampling biases has remained a challenge so far. In this paper, we propose DeepMaxent, which harnesses neural networks to automatically learn shared features among species, using the maximum entropy principle. To do so, it employs a normalised Poisson loss where for each species, presence probabilities across sites are modelled by a neural network. We evaluate DeepMaxent on a benchmark dataset known for its spatial sampling biases, using PO data for calibration and presence-absence (PA) data for validation across six regions with different biological groups and covariates. Our results indicate that DeepMaxent performs better than Maxent and other leading SDMs across all regions and taxonomic groups. The method performs particularly well in regions of uneven sampling, demonstrating substantial potential to increase SDM performances. The method opens the possibility to learn more robust features predicting simultaneously many species to arbitrary large datasets without increased memory requirements. The model likelihood, arising from a Poisson process, makes the method compatible with the integration of more standardised types of data to further increase sampling bias correction. In particular, our approach yields more accurate predictions than traditional single-species models, which opens up new possibilities for methodological enhancement.

**Keywords** species distribution modelling · neural networks · maximum entropy principle · deepmaxent · presence-only data · sampling bias · target-group background

## 1 Introduction

In recent years, the rapid growing number of citizen science projects has contributed significantly to the expansion of biodiversity databases, particularly through the collection of presence-only (PO) observations [Callaghan et al., 2022, Bonnet et al., 2020]. PO records have been instrumental to improve our understanding of species distributions and help to inform conservation strategies Carvalho et al. [2011], Guisan et al. [2013]. However, building reliable Species Distribution Models (SDM) from PO data obtained from opportunistic observations, without a homogeneous sampling

protocol, is challenging due to the many sampling biases and errors to handle (Boakes et al. [2010], Bird et al. [2014], Hughes et al. [2021]). Sampling bias correction in SDM becomes even more challenging when SDM become more complex, e.g. based on deep learning [Deneu et al., 2021, Estopinan et al., 2022, Cole et al., 2023]. The correction of sampling bias in SDMs based on deep learning has not been addressed to date, yet is an essential step if these methods are to take advantage of PO data.

Maxent [Phillips et al., 2006] is one of the most widely used and effective methods for modelling species distributions based on PO data [Warren and Seifert, 2011, Elith et al., 2006, 2020, Valavi et al., 2022]. Maxent generates a probability distribution across sites representing the sampled area (i.e. background points). The probability distribution is defined as a function of various transformations, hereafter called features, of input environmental variables provided by the user. Maxent’s name arises from the fact that it estimates the probability function by maximising its entropy under constraints on the features. A major issue when calibrating any SDM on PO data is spatial sampling bias, which arises from the clustering of PO records in certain areas, typically of higher accessibility or greater human activity. Such bias can distort SDM outputs, leading to inaccurate species distribution estimates [Yackulic et al., 2013, Fithian et al., 2015, Phillips et al., 2009]. To mitigate spatial sampling bias in Maxent, instead of drawing background points uniformly across the area, Phillips et al. [2009] proposed to restrict them to the areas actually sampled, which should reduce the number of false-absences. Specifically, they showed that using occurrences of a group of species, targeted for sampling along with the focal one, to define background points yielded an efficient bias correction (the Target-Group Background correction, hereafter TGB). Furthermore, Maxent is equivalent to a Poisson regression over sites and to inhomogeneous Poisson Point Process (PPP) models [Renner and Warton, 2013], so that common PPP models of the biased sampling of PO data [Fithian et al., 2015] can be used to show that the TGB correction has theoretical guarantees of robustness to spatial sampling bias under some assumptions [Botella et al., 2020].

Deep learning is a family of data-driven methods that allow fitting arbitrary non-linear functions using neural networks, backpropagation and stochastic gradient descent [Goodfellow et al., 2016, Hornik et al., 1989, LeCun et al., 1989]. Unlike other machine learning approaches, deep learning methods allow to simultaneously learn predictive functions while learning an appropriate feature representation, eliminating the need for feature design. Feature design is an important step in traditional SDMs including Maxent [Phillips and Dudík, 2008, Komori et al., 2024], while deep learning-based methods can automatically capture potentially complex and non-linear representations directly from the data [LeCun et al., 2015, Schmidhuber, 2015, Deneu et al., 2021, Estopinan et al., 2022]. Learnt features tend to be more predictive and robust the more species are simultaneously used for training [Chen et al., 2017, Botella et al., 2018], leading to a recent interest in deep learning algorithms for multi-species distribution modelling [Kellenberger et al., 2024]. Deep learning approaches have been shown to have an edge over other methods in cases where structured and high-dimensional data, such as remote sensing imagery, are used as input variables [Deneu et al., 2021, Estopinan et al., 2024]. However, these models are still found to underperform with low-dimensional input variables. After all, they remain susceptible to sampling biases [Zbinden et al., 2024], potentially amplifying them due to their capacity of fitting arbitrary functions, which can lead to inaccurate conclusions and flawed conservation strategies.

In this study, we propose a method that uses the Maxent principle of maximum entropy, with its bias correction capabilities, within a deep learning framework (DeepMaxent) where the features are learnt to predict simultaneously multiple species from PO data. Each species probability distribution across sites is represented as a log-linear function of the joint latent features. Using the equivalence of Maxent and PPPs, we propose a loss function based on the maximum entropy principle of Maxent. We show that computing the loss over any batch of sites preserves the global minimiser of the full loss, so that we can apply the SGD algorithm to train the model and benefit from its regularisation and scalability. We show how the TGB correction can be implemented in DeepMaxent to mitigate spatial sampling bias. We evaluate DeepMaxent on a reference dataset [Elith et al., 2020] encompassing PO data for training and presence-absence (PA) data for evaluation in six distinct regions and comprising different biological groups. We compare the approach to various state-of-art SDMs, notably Maxent, its TGB correction and other multi-species deep learning based SDM. We also test the robustness of DeepMaxent to various neural network architectures and hyper-parameters for its training.

## 2 Materials and Methods

### 2.1 DeepMaxent: point process for SDM based on neural networks

#### 2.1.1 Point processes and Poisson loss

Given some geographic domain  $\mathcal{D} \in \mathbb{R}^2$ , the true abundance of a species  $j \in [1, N]$  across this domain is often modelled by the intensity function  $\lambda_j : \mathcal{D} \mapsto \mathbb{R}_+$  of an inhomogeneous Poisson process [Renner et al., 2015, Fithian et al., 2015]. In this probabilistic model, the expected abundance of species  $j$  in any sub-domain  $d_i \subset \mathcal{D}$  is written  $\lambda_{ij} = \int_{d_i} \lambda_j(z) dz$ . We will generally consider a set of  $K$  non-overlapping sub-domains that we will refer to via their

index  $i \in [1, \dots, K]$ , with  $\mathcal{D} = \bigcup_i d_i$ . We will denote as  $\lambda : \mathcal{D} \mapsto \mathbb{R}_+^N$  the vector-valued function of intensities associated with each point in  $\mathcal{D}$  for a set of  $N$  considered species.

PO data, the actual observations, consist of a set of count vectors  $y := \{y_{ij}\}_{i=1, j=1}^{K, N}$  corresponding to the  $K$  sites and  $N$  species, with  $y_{ij} \in \mathbb{N}_+$ . Most often in PO data, a given  $y_{ij}$  is very sparse.

Since the observation process is modelled as an inhomogeneous Poisson point process, a loss function  $\mathcal{L}$  is defined to integrate the likelihood of the observed counts as a function of the intensity function  $\lambda$  [Cressie, 1993]. In this paper, we use the notation  $\mathcal{L}_{\mathcal{P}}(\lambda, y)$  to denote the Poisson loss between occurrence counts  $y := \{y_{ij}\}_{i=1, j=1}^{K, N}$  and estimated intensity values  $\lambda := \{\lambda_{ij}\}_{i=1, j=1}^{K, N}$  across the  $K$  sites and  $N$  species. This can be formulated as:

$$\mathcal{L}_{\mathcal{P}}(\lambda, y) = \frac{1}{KN} \sum_{i=1}^K \sum_{j=1}^N (\lambda_{ij} - y_{ij} \log \lambda_{ij}). \quad (1)$$

The probability function  $\tilde{\lambda}_j : \mathcal{D} \mapsto \mathbb{R}_+$  of species  $j$ , where  $\int_{\mathcal{D}} \tilde{\lambda}_j(z) dz = 1$ , which describes its relative distribution across space, is defined by normalising  $\lambda_j$  by the partition function over the geographic domain,  $\int_{\mathcal{D}} \lambda_j(z) dz$ . When using the set of sites, the normalised intensity  $\tilde{\lambda}_{ij}$  at site  $i$  for a given species  $j$  is expressed as follows:

$$\tilde{\lambda}_{ij} = \frac{\lambda_{ij}}{\sum_{k=1}^K \lambda_{kj}}. \quad (2)$$

In this context,  $\tilde{\lambda}_{ij}$  can be interpreted as the relative probability of species  $j$  being observed at site  $i$ , given that an occurrence happens somewhere within the  $K$  sites. Equivalently for observation counts, the probability function  $\tilde{y}_{ij}$  for species  $j$ , which represents its relative distribution in space, is defined by normalising  $y_{ij}$  by the sum of observations  $y_{ij}$  over all sites  $K$  in the geographical domain  $\mathcal{D}$ . This normalised intensity  $\tilde{y}_{ij}$  for species  $j$  at site  $i$  is given by :

$$\tilde{y}_{ij} = \frac{y_{ij}}{\sum_{k=1}^K y_{kj}}. \quad (3)$$

When incorporating this normalisation into the Poisson loss, the original formulation of the Poisson log-likelihood can be reformulated to reflect the normalised intensities. In this context,  $\tilde{\lambda}_{ij}$  can be interpreted as the relative probability of species  $j$  being observed at site  $i$ , given that an occurrence happens somewhere within the  $K$  sites. Given that  $\tilde{\lambda}_{ij}$  now represents a proportion (i.e., a "density-like" measure), the loss function is adjusted to ensure that these normalised values are used consistently. The loss can then be expressed as follows:

$$\mathcal{L}_{\mathcal{H}}(\tilde{\lambda}, \tilde{y}) = -\frac{1}{KN} \sum_{i=1}^K \sum_{j=1}^N \left( \frac{y_{ij}}{\sum_{k=1}^K y_{kj}} \right) \log \left( \frac{\lambda_{ij}}{\sum_{k=1}^K \lambda_{kj}} \right). \quad (4)$$

Note that  $\mathcal{L}_{\mathcal{H}}(\tilde{\lambda}, \tilde{y})$  represents a cross entropy loss that corresponds to the Maxent loss formulation, which is also equivalent to the conditional negative log-likelihood of a Poisson point process, as previously demonstrated by Renner and Warton [2013]. As shown in appendix A.1, it implies that the losses in equations 4 and 1 are equivalent in the space of normalised intensities (i.e. values of  $\tilde{\lambda}$ ). However, the Poisson loss introduces a constant factor in the intensity that can affect the learning of the parameters determining  $\tilde{\lambda}$ . In the multi-species case, however, the per-species normalisation of the PO labels  $y_{ij}$  according to the total number of observations  $\sum_{k=1}^K y_{kj}$  leads to all species having the same weight in the total loss. This can lead to learning a joint latent representation (see below) that is overwhelmed by the noisy learning signal from rare species. To address this problem, we propose to weight the loss per species proportionally to their number of PO occurrences  $w_j := \sum_{i=1}^K y_{ij}$ , which leads to:

$$\mathcal{L}_{\mathcal{H}, \mathcal{W}}(\tilde{\lambda}, y) = -\frac{1}{KN} \sum_{i=1}^K \sum_{j=1}^N w_j \left( \frac{y_{ij}}{\sum_{k=1}^K y_{kj}} \right) \log \left( \frac{\lambda_{ij}}{\sum_{k=1}^K \lambda_{kj}} \right) = -\frac{1}{KN} \sum_{i=1}^K \sum_{j=1}^N y_{ij} \log \left( \frac{\lambda_{ij}}{\sum_{k=1}^K \lambda_{kj}} \right) \quad (5)$$

Note that if the intensity of each species was based on pre-determined covariates, this weighting wouldn't have any impact on the estimated intensities. However, given that we learn a joint features representation between species (as

introduced below), the model can calibrate this representation to favour performance on certain species at the expense of others during learning.

### 2.1.2 Feature extraction based on neural network

In Maxent, the intensity value in a given site is a function of a vector of environmental variables  $x \in \mathbb{R}^P$  in that site. More precisely, this intensity is defined as a log-linear function of a feature vector  $f(x)$ , composed of pre-determined transformations of  $x$ . We extend the principle of Maxent by replacing its feature vector  $f(x)$  with a feature extractor instantiated as a neural network  $g : \mathbb{R}^P \mapsto \mathbb{R}^C$  parametrised by  $\theta$ , where  $C$  is the dimensionality of the last hidden layer, the joint latent representation that is common to all species, as illustrated in Figure 1. The intensity function  $\lambda_j(x)$  of species  $j$  in DeepMaxent is then given by:

$$\lambda_j(x) = \exp \left( \sum_{c=1}^C \gamma_{jc} g(x; \theta)_c + b_j \right), \quad (6)$$

where  $\gamma \in \mathbb{R}^{N \times C}$  and  $b \in \mathbb{R}^N$  can be treated the weights and the bias term of an additional linear layer that maps from  $\mathbb{R}^C$  to  $\mathbb{R}^N$ , corresponding to the number of species. The function  $g$  can automatically learn complex, non-linear relationships between environmental variables and the presence of multiple species from the data, potentially enabling the model to identify environmental patterns not considered by traditional approaches. A key advantage of DeepMaxent is its scalability to many species. DeepMaxent computes a single feature vector  $g(x)$  per site for all species, so adding more species does not increase the computational cost of computing  $g(x)$ . In addition, all the species contribute to learning  $g$ , which means that species with very few observations can benefit from the learning signal provided by other species. These properties are interesting when  $g$  is designed to have more capacity, e.g. when it includes several hidden layers or if the neural network takes high dimensional input variables such as spatio-temporal data.

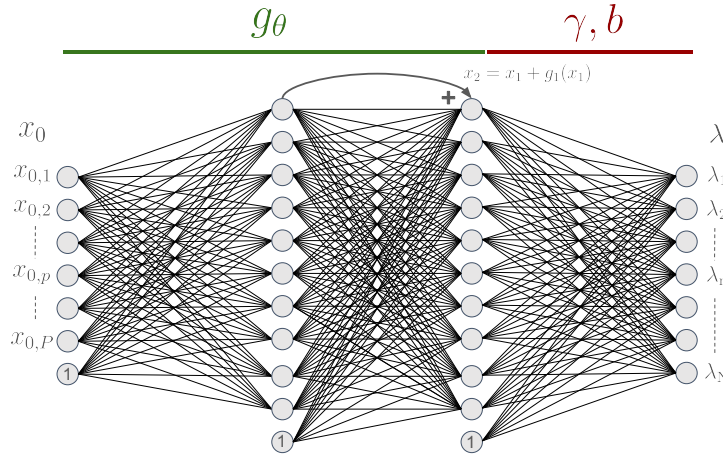


Figure 1: The residual neural network to estimate the intensity  $\lambda$  from variable input ( $x_0$ ), where  $P$  as the variable number input, where  $C$  is the number of hidden layer nodes, and  $N$  denotes the number of species (or target categories). The illustrated case involves two hidden layers. In the special case where there is only one hidden layer, no residual addition is applied.

DeepMaxent is architecture-agnostic, and in this work we chose a  $g$  based on a residual neural network architecture (see Figure 1), which takes the vector of environmental variables  $x$  as input.  $x$  is transformed by a sequence of hidden layers  $g = g_1 \circ g_2 \circ \dots \circ g_L$ , whose number  $L$  and size are the main hyper-parameters. Residual connections between adjacent hidden layers were incorporated introducing shortcuts that allow the input of a layer to bypass the non-linear transformations and be directly added to the output ( $x_{l+1} = x_l + g_l(x_l)$ ). The output of the last hidden layer ( $g(x)$ ), used as our latent feature vector, leads to the output layer, which is composed of each species intensity measure  $\lambda_j(x)$ .

### 2.1.3 Batched algorithm and partition function approximation in DeepMaxent

One of the challenges of adapting Maxent to a deep learning framework is the computation of the partition functions  $\sum_k \lambda_{k,j}$  for each species  $j$ , corresponding to the denominator in equation 2, which normalises the predicted intensity



$\lambda_j(x)$  over the  $K$  sites. When the geographic domain becomes large, or when the spatial resolution is increased, the number of terms  $K$  becomes large, and computing this function may be particularly expensive in a deep learning context where each  $g_\theta(x_i)$  may already be costly. To address this consistently with the batched optimisation algorithms used in deep learning, such as Stochastic Gradient Descent (SGD), we modify the normalisation step by computing the loss function on a small random subset of sites  $B \subset [1, K]$ , called batch, hence normalising the intensities within the batch. The batch-wise loss is then written as follows:

$$\mathcal{L}(\tilde{\lambda}_{i \in B}, \tilde{y}_{i \in B}) = -\frac{1}{|B|N} \sum_{i \in B} \sum_{j=1}^N \frac{y_{ij}}{\sum_{i \in B} y_{ij}} \log \left( \frac{\lambda_j(x_i)}{\sum_{i \in B} \lambda_j(x_i)} \right). \quad (7)$$

This batch-wise loss makes the model optimisation computationally feasible by computing it one batch at a time, and thus allows it to be scalable to large domains and trained efficiently with any stochastic optimiser based on random batches, such as the very commonly used method Adam which we use here (Kingma [2014]). However, it must be noted that the batch-wise loss is not a simple approximation of the full loss, as for instance the normalised labels tend to have a larger value for smaller batches. Nevertheless, we provide the mathematical guarantee that, for any batch size  $n$  ( $1 < n < K$ ), optimising the model on all batches also optimises the full loss (see Appendix A.3). This suggests that our final estimator should be somehow close to the global minimiser of the full loss, even though we will unlikely obtain the latter. Similarly to supervised contrastive methods Khosla et al. [2020], the intensity predictions could become more specific, and more concentrated around the occurrences, as the batch size increases. A small batch size could, on the other hand, result in smoother species intensities. Therefore we tested the impact of the batch size during training on the final predictions of DeepMaxent. Given the expression of  $\lambda_j(x_i)$  in equation 6, its normalisation across the batch  $B$  writes as  $\exp(\gamma_j^T g_\theta(x_i)) / \sum_{k \in B} \exp(\gamma_j^T g_\theta(x_k))$ . This is a particular case of a softmax function applied to the logits  $\gamma_j^T g_\theta(x_k)$  over the batch  $B$ .

#### 2.1.4 Spatial sampling bias correction with Target-Group Background correction

When occurrence concentration is biased by spatial variations in sampling effort, a popular SDM correction approach is the Target-Group Background (TGB) method Phillips et al. [2009], which was initially proposed to correct sampling bias in Maxent. The method basically approximates the spatial sampling effort through the distribution of occurrences of a Target-Group of species, providing background points to Maxent for each site where TG species were reported. In other words, TGB restricts the study domain to the sites with at least an evidence of sampling effort (one observation), which reduces the problem of false absences. The strategy is expected to work when TG species are reported jointly with the focal species (e.g. the TG is a biological group targeted by a citizen science program).

In DeepMaxent, which is a multi-species model, it can be relevant to define the whole set of  $N$  predicted species as the Target-Group. This is done in the empirical study below: For each region, each biological group is the set of predicted species and the Target-Group, so the sites where at least one of the these TG species was observed are the background sites. Then, for each iteration of our batched algorithm, we randomly select a same subset of TG background sites for all predicted species.

#### 2.1.5 L2-regularisation implementation in DeepMaxent

Maxent makes use of a L1 penalisation, on the species weights  $\gamma_j$  that model the relation between the features and the density prediction. The L1 term, known as LASSO penalty, encourages  $\gamma_j$  to become sparse, thus selecting a subset of features.

The L1 regularisation is important in Maxent due to a number of features that grows more than quadratically with the number of environmental variables provided a priori by the user. In DeepMaxent, the latent features are learned to maximise prediction performances, removing the need for feature selection. For DeepMaxent, we employ a L2 regularisation, i.e. a penalty on the Euclidean norm of the  $\gamma_j$ , which encourages small but non-zero weights, which may induce a smoothing of the estimated species intensities. The total loss function thus becomes:

$$\mathcal{L}_{\text{total}}(\tilde{\lambda}, y; \theta, \gamma) = \mathcal{L}_{\mathcal{H}}(\tilde{\lambda}, y; \theta, \gamma) + \frac{\tau}{2} (\|\theta\|_2^2 + \|\gamma\|_2^2) \quad (8)$$

where  $\tau$  is the weight decay coefficient. This term penalises large weight values, encouraging the model to learn smaller weights without enforcing sparsity.

## 2.2 Evaluation of model performance

### 2.2.1 Dataset

For our experiments we used an openly released dataset from the National Centre for Ecological Analysis and Synthesis (NCEAS) [Elith et al., 2020]. This dataset includes PO (for SDM training) and presence-absence (PA, for SDM evaluation) data from six global regions: Australian Wet Tropics (AWT), Canada (CAN), New South Wales (NSW), New Zealand (NZ), South America (SA) and Switzerland (SWI) [Elith et al., 2020]. Each regions is associated with a specific set of species, and sometimes from several biological groups (see Table 1), with a total of 226 anonymous species. The dataset provides specific environmental variables for each region, including climatic, soil or location variables (see more details Elith et al. [2020]).

The NCEAS dataset has been used to evaluate and compare various SDM methods based on presence-only data [Elith et al., 2006, Phillips et al., 2009, Zbinden et al., 2024, Valavi et al., 2022] and it allows us to compare the performances of DeepMaxent with existing SDM methods (section below).

[Phillips et al., 2009] studied spatial sampling biases and found that the PO data in certain regions (AWT, CAN, and SWI) contain high levels of such biases. This is the data that we use to train all our models.

Table 1: The total number of species, the occurrence number in PO data and the total number of species presence in PA data for each region and biological group

Code	Location	Biological Group	Species number	Occurrences number	
				PO	PA
AWT	Australian wet tropics	bird	20	3105	340
AWT	Australian wet tropics	plant	20	701	102
CAN	Ontario, Canada	bird	20	5063	14571
NSW	New South Wales	bates	10	187	570
NSW	New South Wales	birds	7	1781	1839
NSW	New South Wales	plants	29	680	5329
NSW	New South Wales	reptile	8	675	1008
NZ	New Zealand	plant	52	3088	19120
SA	South America	plant	30	2220	152
SWI	Switzerland	tree	30	35105	10013

### 2.2.2 Evaluation metrics

To directly compare our results to Phillips et al. [2009], Zbinden et al. [2024] and Valavi et al. [2022], we evaluated our method performances with the Area Under the ROC Curve (AUC) on the PA plots of the NCEAS dataset. The AUC is the empirical probability that a presence site has a higher model-predicted value than an absence site. In other words, it measures the model ability to distinguish between presence and absence classes based on its predicted scores. The NCEAS dataset being decomposed into regions and biological groups, for each region the average across biological groups of the AUCs and the general average of the latters across regions were calculated.

### 2.2.3 Implementation details

In this study, the first version of DeepMaxent method was implemented using Python and PyTorch framework and openly available at <https://github.com/RYCKEWAERT/deepmaxent>. The architecture was designed as a Multilayer Perceptron (MLP) residual neural networks (see Figure 1) with four fully connected layers interconnected by Rectified Linear Unit (ReLU) activation.

Cross-validation was performed to optimise hyperparameters, using spatial blocking based on geographic data [Valavi et al., 2019, Roberts et al., 2017]. As suggested in Zbinden et al. [2024], the cross-validation is performed exclusively on PO data. Once cross-validation has been performed, the parameter values are set and are the same for all DeepMaxent models for all regions and target-Groups, so as to be comparable with other models. The final model is then calibrated with these hyperparameter values on the whole PO data and applied to the PA data.

### 2.2.4 Baseline losses

We implemented the Poisson regression loss (see Eq. 1) to test the effect of the density normalisation in DeepMaxent on the estimator quality. The Poisson loss is also equivalent to fitting the intensity of a gridded inhomogeneous Poisson

point process (Renner and Warton [2013]) over our sites. Other commonly used loss functions in deep learning, and notably for SDMs, namely Cross-Entropy over species (CE, Deneu et al. [2021], Brun et al. [2024]) and Binary Cross-Entropy (BCE, Benkendorf and Hawkins [2020], Zbinden et al. [2024]) were implemented. All the baseline implementations were performed using the same architecture (2 hidden layers), optimiser (Adam), and hyperparameters (learning rate: 0.0002, batch size: 250) as the best performing DeepMaxent implementation, except for the weight decay  $\tau$  which we kept to 0 for the baselines, as this regularisation did not improve any of them. Each loss function was evaluated both with and without TGB correction.

The Cross-Entropy over species loss (CE),  $\mathcal{L}_{\text{CE}}(\lambda, y)$  (equation 9), measures for each site the deviation between a predicted probability distribution across species and the associated empirical distribution based on the species observations in that site. In this case, the predicted probabilities are obtained by normalising the intensity values over the species  $\lambda := \{\lambda_{ij}\}_{i=1, j=1}^{K, N}$  (implemented with a softmax as for DeepMaxent):

$$\mathcal{L}_{\text{CE}}(\lambda, y) = -\frac{1}{K} \sum_{i=1}^K \sum_{j=1}^N \frac{y_{ij}}{\sum_{k=1}^N y_{ik}} \log \left( \frac{\lambda_{ij}}{\sum_{k=1}^N \lambda_{ik}} \right) \quad (9)$$

The Binary Cross-Entropy loss,  $\mathcal{L}_{\text{BCE}}(\lambda, y)$ , was implemented for the case where  $y \in \{0, 1\}$  represents a binary variable, taking the value 1 for presence and 0 for absence. For binary classification, a sigmoid function is used to map the predicted logits  $\lambda_{ij}$  to probabilities, making it a simpler case compared to multi-class classification (which uses the softmax function). The BCE loss is defined as:

$$\mathcal{L}_{\text{BCE}}(\lambda, y) = -\frac{1}{KN} \sum_{i=1}^K \sum_{j=1}^N (y_{ij} \log \sigma(\log(\lambda_{ij})) + (1 - y_{ij}) \log(1 - \sigma(\log(\lambda_{ij})))), \quad (10)$$

where  $\sigma$  is the sigmoid function here applied to the linear predictor, or "logit", of each species  $\log(\lambda_{ij}) = \gamma_j^T g_{\theta}(x_i) + b_j$ .

### 3 Results

#### 3.1 Comparative analysis of SDM methods

Table 2 shows the performances of various standard comparative methods, including Maxent, Boosted Regression Tree (BRT) with or without TGB correction, and the neural network model for multi-species proposed by Zbinden et al. [2024], all evaluated with the average AUC per region and the general average [Phillips et al., 2006, Valavi et al., 2022, Zbinden et al., 2024]. It also contains the performance of our main DeepMaxent implementation (best architecture and hyperparameters) and the baseline losses, with or without the TGB correction.

Without TGB sampling bias correction, performances are close among methods and range from 0.718 to 0.723 in general average AUC (Table 2), except for the CE loss which is much better (0.731). Except for the latter, we observe no general performance gain for the tested deep learning losses (BCE, Poisson, DeepMaxent, ranging from 0.719 to 0.720) compared to the literature methods, e.g. Maxent (0.721) or the best SDM Ensemble of Valavi et al. [2022] (0.723).

The TGB correction brings a consistent general performance improvement for all methods, including the ones of the literature and our implementations. However, not all approaches respond equally strongly to the sampling bias correction. For instance, Maxent gains 0.039 in general averaged AUC by using TGB, and it is the same for BRT. Regarding our implemented baseline losses, TGB induces an AUC gain of 0.011 for the CE loss, 0.044 for the BCE loss and 0.039 for the Poisson loss. Finally, DeepMaxent gains 0.047, making it the best method in terms of general AUC (0.767). These results show that the proposed DeepMaxent is particularly adapted to this bias correction technique while it enables to leverage the predictive potential of neural networks for multi-species spatial intensity estimation. The largest region average AUC gains were mostly seen in regions CAN and AWT, where spatial sampling bias is the strongest according to [Phillips et al., 2009]. Note that the best method of Zbinden et al. [2024], achieving a general averaged AUC of 0.755, incorporated both random and TGB points as absences in their BCE loss, and their results specifically showed the key role of the TGB points in this performance. DeepMaxent-TGB also had the best AUC in four of the six regions (CAN, NSW, NZ, SWI), showing that it is robust across regions and biological groups (NSW includes four biological groups, see Table 1). BCE-TGB is the second-best method in general AUC (0.763). In contrast, CE-TGB and Poisson-TGB yield poorer general AUC (0.745 and 0.758) than Maxent-TGB (0.760) or BRT-TGB (0.759).

Table 2: Comparison of method performance by region-averaged AUC and general averaged AUC over all regions. The best average AUC for each column is highlighted in bold, while the second-best averaged AUC is underlined. The references correspond to results from the following articles: [1] Valavi et al. [2022], [2] Phillips et al. [2009] and [3] Zbinden et al. [2024].

	Regions						avg
	AWT	CAN	NSW	NZ	SA	SWI	
<b>Results from the literature</b>							
<i>Single-species models</i>							
Maxent [1]	0.686	0.587	0.700	0.738	0.804	0.809	0.721
BRT [1]	0.681	<u>0.577</u>	0.701	0.735	0.795	0.816	0.718
RF down-sampled [1]	0.675	0.572	0.715	0.746	<u>0.813</u>	0.818	0.723
Ensemble [1]	0.683	0.580	0.710	<u>0.749</u>	0.806	0.812	0.723
Maxent (using TGB) [2]	<b>0.732</b>	0.716	0.741	0.738	0.798	0.837	0.760
BRT (using TGB) [2]	0.700	<u>0.728</u>	0.738	0.740	0.792	0.842	0.757
<i>Multi-species models</i>							
Zbinden et al. [3]	0.704	0.714	0.719	0.741	<b>0.815</b>	0.838	0.755
<b>Results from our implementations</b>							
<i>Baseline losses</i>							
CE	0.701	0.661	0.732	0.724	0.772	0.793	0.731
CE (using TGB)	<u>0.727</u>	0.708	0.739	0.732	0.771	0.792	0.745
BCE	0.656	0.600	0.718	0.736	0.804	0.799	0.719
BCE (using TGB)	0.723	0.726	<u>0.743</u>	0.739	0.803	<u>0.846</u>	<u>0.763</u>
Poisson loss	0.658	0.599	0.714	0.737	0.804	0.799	0.719
Poisson loss (using TGB)	0.712	0.727	0.732	0.731	0.800	<u>0.846</u>	0.758
<i>Proposed loss</i>							
DeepMaxent	0.654	0.593	0.718	0.744	0.803	0.810	0.720
<b>DeepMaxent (using TGB)</b>	0.714	<b>0.732</b>	<b>0.752</b>	<b>0.754</b>	0.803	<b>0.850</b>	<b>0.767</b>

### 3.2 Sensitivity study

Table 3 shows the average AUCs calculated for all regions, according to six different values for each hyperparameter: batch size, number of hidden layers and weight decay. A detailed analysis of AUC values for each region is provided in the appendix (D). The general performance of DeepMaxent-TGB was quite robust to hyper-parameter choices and neural network depth. In particular, DeepMaxent-TGB kept a general average AUC above 0.764, i.e. above all other methods, for all tested batch sizes, ranging from 10 to 2500. Qualitatively, a smaller batch size induces smoother species intensity maps, while larger batch size tends to concentrate the intensity in higher abundance areas, as illustrated for one species in the region CAN in Figure 2. The L2 regularisation (weight decay) has a larger impact on DeepMaxent-TGB performance. The regularisation has a small but consistently positive impact on the performance up to a value of  $3 \times 10^{-4}$  (see Table 3), while further increasing the weight decay value results on a performance degradation due to oversmoothing (see Figure 2). DeepMaxent-TGB keeps the highest general average AUC among tested methods for  $\tau$  ranging from  $3 \times 10^{-4}$  to  $1 \times 10^{-3}$ . Varying the number of hidden layers in the neural network architecture of DeepMaxent from one to two had almost no effect, with a same general averaged AUC of 0.767 (Table 3), and the AUC softly and progressively decreased for three (0.766), four (0.764), five (0.762) and six layers (0.759).

Table 3: Average AUC values for DeepMaxent-TGB across all regions, calculated for six different values of each hyperparameter: batch size, number of hidden layers, and weight decay. The default values used are a batch size of 250, two hidden layers, and a weight decay of  $3e-4$ .

Batch size	AUC	Hidden layers	AUC	Weight decay	AUC
10	0.765	1	0.767	0	0.762
25	0.765	2	0.767	3e-5	0.763
100	0.767	3	0.766	1e-4	0.765
250	0.767	4	0.764	3e-4	0.767
1000	0.766	5	0.762	1e-3	0.765
2500	0.764	6	0.759	3e-3	0.757

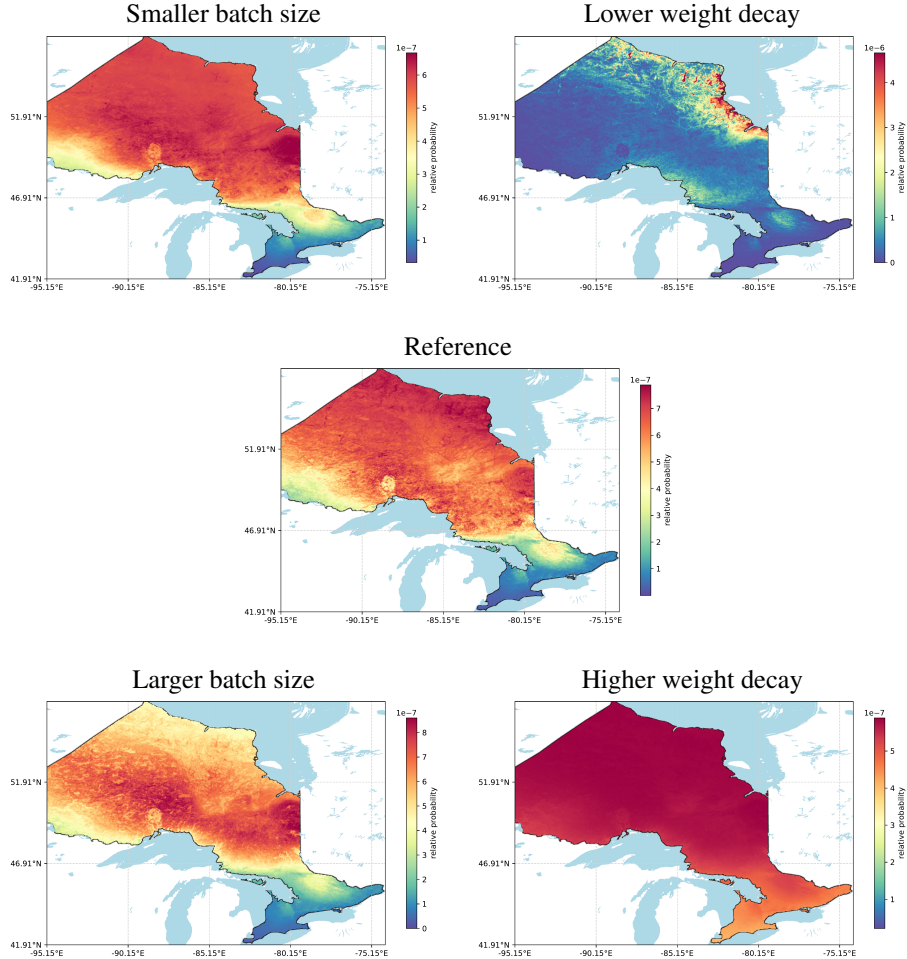


Figure 2: Estimated probabilities for the species can14 (CAN) by varying batch size and weight decay, while keeping other hyperparameters at their default values (batch size = 250, hidden layers = 2, weight decay =  $3 \times 10^{-4}$ ) corresponding to an AUC of 0.919. Results are shown (i) for different batch sizes: (i) a lower batch size (10) with AUC of 0.923, (ii) a reduced weight decay ( $1 \times 10^{-5}$ ) with an AUC of 0.904, and (iii) a higher batch size (2500) with an AUC of 0.921, while a larger weight decay ( $1 \times 10^{-1}$ ) with an AUC of 0.555.

## 4 Discussion

The proposed DeepMaxent method, which jointly learns the spatial density for multiple species through a scalable batched algorithm and allows inducing results with maximum entropy, was evaluated on a benchmark consisting of six distinct regional datasets. These include different biological groups and covariates. Compared to state-of-the-art methods, DeepMaxent with the TGB bias correction achieved the highest average AUC across all regions and outperforms all other methods in four of the six regions (see Table 2). It notably surpasses the single-species Maxent and BRT methods with TGB correction, as well as the recent multi-species neural network method of Zbinden et al. [2024]. We applied the TGB correction for spatial sampling bias correction with presence-only data Phillips et al. [2009] to DeepMaxent by restricting training sites to the ones including at least one observation of the biological group from which the species are predicted. We showed that DeepMaxent is particularly adapted to the TGB correction as it outperformed the same multi-species model learned with three other commonly used loss functions in this context: The Poisson loss from which it derives (Renner and Warton [2013]), the cross-entropy over species (Deneu et al. [2021], Brun et al. [2024]) and the binary-cross-entropy loss (encoding occurrences as presence and pseudo-absences, Benkendorf and Hawkins [2020], Zbinden et al. [2024]). We believe this is an important result in the sense that it demonstrates that (i) our methodology for extending Maxent to deep learning is functional on biased presence-only data, (ii) it outperforms the original Maxent method on the same input tabular data and (iii) it outperforms alternative

loss functions used for multi-species SDM based on neural networks. Further efforts could adapt different background sampling strategies while respecting dataset specificity for biological groups and species [Schartel and Cao, 2024].

For feature extraction, a fully-connected neural network with just one or two hidden layers was found to be optimal for these datasets which comprise only tabular data. This simple MLP architecture performs well on tasks involving low-dimensional and unstructured data. While deep learning is often associated with large models, simpler architectures can be more efficient and beneficial for such tasks. For multi-species settings, DeepMaxent reduces the overall computational cost by learning a single feature extractor shared across species [Ba and Caruana, 2014, Raghu et al., 2017], notably compared to methods relying on more resource-intensive, per-species, calculations [Merow et al., 2014].

More broadly, the proposed approach is flexible regarding the type of input and species observation data and should facilitate data integration approaches in the future by using neural networks. In the case of more structured input data, the neural network architecture could be adapted to ingest different types of structured inputs (e.g., sequential or spatial), which might lead to capturing complementary spatio-temporal environmental patterns (Deneu et al. [2021], Estopinan et al. [2022]). Importantly, DeepMaxent is the first method to bridge the gap between deep learning and point process-based SDM (Renner et al. [2015]). Using the formulation of the Poisson process within the loss function, our method could seamlessly incorporate additional loss terms to integrate diverse and more standardised types of species observations, such as presence-absence surveys (Fithian et al. [2015]), detection/non-detection histories (Koshkina et al. [2017]), and abundance or imperfect count data (Dorazio [2014]). This approach, called integrated SDMs, have recently been highlighted as a promising avenue to enhance the reliability of SDMs (Miller et al. [2019], Isaac et al. [2020], Mostert and O’Hara [2023]). Extending DeepMaxent with this approach could use standardised datasets to disentangle the real abundance of each species from detection biases, while harnessing the extensive geographic coverage of opportunistic presence-only data.

The DeepMaxent loss was notably more adapted to the TGB bias correction technique compared to the BCE, CE and Poisson losses. Regarding BCE with TGB, which was the second-best method overall, it suggests that interpreting TG background sites as absences is reasonable as it filters many false absences, but not optimal as it gives too much importance to the remaining false absences. An apparent paradox is that, given the equivalence between the Maxent and Poisson process losses (Renner and Warton [2013]), which applies to any model for the spatial intensity (e.g. neural networks) and induces the equivalence between the DeepMaxent-TGB and Poisson-TGB losses, we would expect similar results for the two latter losses, but they were different. We see two hypothetical algorithmic explanations that could jointly play a role in this discrepancy: (i) the intensity scale in the Poisson loss, which cancels in the DeepMaxent loss, affects the gradient on the other parameters determining the density, the learning trajectory and final estimates, and (ii) the stochastic batched algorithm interacts with (i) in deviating the two estimates. Even though the CE loss had the best results without TGB correction, and thus appears natively less sensitive to spatial sampling bias, its performances with TGB remained below the best methods which relied on this correction (DeepMaxent, BCE, Poisson, Maxent). This limited performance may be due to the loss of information on the spatial variations of the intensity for each species when normalising the intensity across species for each site. More broadly, regarding the estimation of multiple species spatial intensities from biased presence-only data, it suggests that learning to classify the most likely observed species per site is a limited approach.

Additionally to the smoothing of species spatial densities brought by the loss function of DeepMaxent, which induces entropy maximisation, this study also highlighted that the L2 regularisation induced more spatial smoothing which enhances model performance in a deep learning setting. This spatial smoothing can be related to the effect of the L2 regularisation of Maxent Phillips et al. [2006] and, more broadly, to regularised loss functions improving the robustness of point process based SDMs Komori et al. [2024]. However, excessively high weight decay values ( $\tau$ ) can be detrimental, resulting penalising model weights too harshly. This can lead to a over-smoothed species spatial densities which don’t capture the actual spatial variations of a species abundance.

We provided a mathematical guarantee to justify the use of a stochastic batched gradient descent algorithm to learn DeepMaxent on batches of sites, and we further showed that varying the batch size (from 10 to 2500) had a relatively weak impact on general performances. From a computational perspective, such algorithm drastically reduces the need for computer memory when training DeepMaxent compared to optimising the loss function on all sites at each iteration, as only the data for one batch needs to be loaded at a time in memory. Yet, the batch size may affect the learning trajectory due to the approximation of the partition function, and we noticed its influence on the final model performance. Similarly to the increasing the weight decay, we observed that decreasing the batch size may smooth species spatial densities. Although batch size is one of the less sensitive hyperparameters, identifying an optimal batch size remains important and it may interact with e.g. the learning rate and number of epochs. The use of the algorithm Adam Kingma [2014], known for its robustness to hyperparameter adjustments, contributed to the stability of DeepMaxent performances.

One classical limitation of PO-based SDMs is the fact that, in many cases, no PA data is available to validate the model hyperparameters. Although we have seen that DeepMaxent is rather robust to hyperparameter choices, including mini-batch size, number of hidden layers, and weight decay value, finding the right hyperparameters is required to obtain optimal results. In order to mitigate this limitation, we verified that hyper-parameter tuning by validating the results on PO data results on good values being chosen, thus removing the need for PA data during model training. These results can be found in the Appendix C.

## References

- Corey T. Callaghan, Thomas Mesaglio, John S. Ascher, Thomas M. Brooks, Analyn A. Cabras, Mark Chandler, William K. Cornwell, Indiana Cristóbal Ríos-Málaver, Even Dankowicz, and Naufal Urfi Dhiya'ulhaq. The benefits of contributing to the citizen science platform iNaturalist as an identifier. *PLoS biology*, 20(11):e3001843, 2022. URL <https://journals.plos.org/plosbiology/article?id=10.1371/journal.pbio.3001843>. Publisher: Public Library of Science San Francisco, CA USA.
- Pierre Bonnet, Alexis Joly, Jean-Michel Faton, Susan Brown, David Kimiti, Benjamin Deneu, Maximilien Servajean, Antoine Affouard, Jean-Christophe Lombardo, Laura Mary, Christel Vignau, and François Munoz. How citizen scientists contribute to monitor protected areas thanks to automatic plant identification tools. *Ecological Solutions and Evidence*, 1(2):e12023, December 2020. ISSN 2688-8319, 2688-8319. doi:10.1002/2688-8319.12023. URL <https://besjournals.onlinelibrary.wiley.com/doi/10.1002/2688-8319.12023>.
- Sílvia B. Carvalho, José C. Brito, Eduardo G. Crespo, Matthew E. Watts, and Hugh P. Possingham. Conservation planning under climate change: Toward accounting for uncertainty in predicted species distributions to increase confidence in conservation investments in space and time. *Biological Conservation*, 144(7):2020–2030, July 2011. ISSN 0006-3207. doi:10.1016/j.biocon.2011.04.024. URL <https://www.sciencedirect.com/science/article/pii/S0006320711001649>.
- Antoine Guisan, Reid Tingley, John B. Baumgartner, Ilona Naujokaitis-Lewis, Patricia R. Sutcliffe, Ayesha I. T. Tulloch, Tracey J. Regan, Lluís Brotons, Eve McDonald-Madden, Chrystal Mantyka-Pringle, Tara G. Martin, Jonathan R. Rhodes, Ramona Maggini, Samantha A. Setterfield, Jane Elith, Mark W. Schwartz, Brendan A. Wintle, Olivier Broennimann, Mike Austin, Simon Ferrier, Michael R. Kearney, Hugh P. Possingham, and Yvonne M. Buckley. Predicting species distributions for conservation decisions. *Ecology Letters*, 16(12):1424–1435, 2013. ISSN 1461-0248. doi:10.1111/ele.12189. URL <https://onlinelibrary.wiley.com/doi/abs/10.1111/ele.12189>. \_eprint: <https://onlinelibrary.wiley.com/doi/pdf/10.1111/ele.12189>.
- Elizabeth H Boakes, Philip JK McGowan, Richard A Fuller, Ding Chang-qing, Natalie E Clark, Kim O'Connor, and Georgina M Mace. Distorted views of biodiversity: spatial and temporal bias in species occurrence data. *PLoS biology*, 8(6):e1000385, 2010.
- Tomas J Bird, Amanda E Bates, Jonathan S Lefcheck, Nicole A Hill, Russell J Thomson, Graham J Edgar, Rick D Stuart-Smith, Simon Wotherspoon, Martin Krkosek, Jemina F Stuart-Smith, et al. Statistical solutions for error and bias in global citizen science datasets. *Biological Conservation*, 173:144–154, 2014.
- Alice C Hughes, Michael C Orr, Keping Ma, Mark J Costello, John Waller, Pieter Provoost, Qinmin Yang, Chaodong Zhu, and Huijie Qiao. Sampling biases shape our view of the natural world. *Ecography*, 44(9):1259–1269, 2021.
- Benjamin Deneu, Maximilien Servajean, Pierre Bonnet, Christophe Botella, François Munoz, and Alexis Joly. Convolutional neural networks improve species distribution modelling by capturing the spatial structure of the environment. *PLOS Computational Biology*, 17(4):e1008856, April 2021. ISSN 1553-7358. doi:10.1371/journal.pcbi.1008856. URL <https://dx.plos.org/10.1371/journal.pcbi.1008856>. Number: 4.
- Joaquim Estopinan, Maximilien Servajean, Pierre Bonnet, François Munoz, and Alexis Joly. Deep species distribution modeling from sentinel-2 image time-series: a global scale analysis on the orchid family. *Frontiers in Plant Science*, 13:839327, 2022.
- Elijah Cole, Grant Van Horn, Christian Lange, Alexander Shepard, Patrick Leary, Pietro Perona, Scott Loarie, and Oisín Mac Aodha. Spatial implicit neural representations for global-scale species mapping. In *International Conference on Machine Learning*, pages 6320–6342. PMLR, 2023.
- Steven J. Phillips, Robert P. Anderson, and Robert E. Schapire. Maximum entropy modeling of species geographic distributions. *Ecological Modelling*, 190(3):231–259, January 2006. ISSN 0304-3800. doi:10.1016/j.ecolmodel.2005.03.026. URL <https://www.sciencedirect.com/science/article/pii/S030438000500267X>.
- Dan L. Warren and Stephanie N. Seifert. Ecological niche modeling in Maxent: the importance of model complexity and the performance of model selection criteria. *Ecological Applications*, 21(2):335–342, 2011. ISSN 1939-5582.

- doi:10.1890/10-1171.1. URL <https://onlinelibrary.wiley.com/doi/abs/10.1890/10-1171.1>. \_eprint: <https://onlinelibrary.wiley.com/doi/pdf/10.1890/10-1171.1>.
- Jane Elith, Catherine H. Graham, Robert P. Anderson, Miroslav Dudík, Simon Ferrier, Antoine Guisan, Robert J. Hijmans, Falk Huettmann, John R. Leathwick, Anthony Lehmann, Jin Li, Lucia G. Lohmann, Bette A. Loiselle, Glenn Manion, Craig Moritz, Miguel Nakamura, Yoshinori Nakazawa, Jacob McC. M. Overton, A. Townsend Peterson, Steven J. Phillips, Karen Richardson, Ricardo Scachetti-Pereira, Robert E. Schapire, Jorge Soberón, Stephen Williams, Mary S. Wisz, and Niklaus E. Zimmermann. Novel methods improve prediction of species' distributions from occurrence data. *Ecography*, 29(2):129–151, April 2006. ISSN 0906-7590, 1600-0587. doi:10.1111/j.2006.0906-7590.04596.x. URL <https://onlinelibrary.wiley.com/doi/10.1111/j.2006.0906-7590.04596.x>.
- Jane Elith, Catherine Graham, Roozbeh Valavi, Meinrad Abegg, Caroline Bruce, Andrew Ford, Antoine Guisan, Robert J. Hijmans, Falk Huettmann, Lucia Lohmann, Bette Loiselle, Craig Moritz, Jake Overton, A. Townsend Peterson, Steven Phillips, Karen Richardson, Stephen Williams, Susan K. Wiser, Thomas Wohlgemuth, and Niklaus E. Zimmermann. Presence-only and Presence-absence Data for Comparing Species Distribution Modeling Methods. *Biodiversity Informatics*, 15(2):69–80, July 2020. ISSN 1546-9735. doi:10.17161/bi.v15i2.13384. URL <https://journals.ku.edu/jbi/article/view/13384>.
- Roozbeh Valavi, Gurutzeta Guillera-Arroita, José J. Lahoz-Monfort, and Jane Elith. Predictive performance of presence-only species distribution models: a benchmark study with reproducible code. *Ecological Monographs*, 92(1):e01486, February 2022. ISSN 0012-9615, 1557-7015. doi:10.1002/ecm.1486. URL <https://esajournals.onlinelibrary.wiley.com/doi/10.1002/ecm.1486>.
- Charles B. Yackulic, Richard Chandler, Elise F. Zipkin, J. Andrew Royle, James D. Nichols, Evan H. Campbell Grant, and Sophie Veran. Presence-only modelling using MAXENT: when can we trust the inferences? *Methods in Ecology and Evolution*, 4(3):236–243, 2013. ISSN 2041-210X. doi:10.1111/2041-210x.12004. URL <https://onlinelibrary.wiley.com/doi/abs/10.1111/2041-210x.12004>. \_eprint: <https://onlinelibrary.wiley.com/doi/pdf/10.1111/2041-210x.12004>.
- William Fithian, Jane Elith, Trevor Hastie, and David A. Keith. Bias correction in species distribution models: pooling survey and collection data for multiple species. *Methods in Ecology and Evolution*, 6(4):424–438, April 2015. ISSN 2041-210X, 2041-210X. doi:10.1111/2041-210X.12242. URL <https://onlinelibrary.wiley.com/doi/10.1111/2041-210X.12242>.
- Steven J. Phillips, Miroslav Dudík, Jane Elith, Catherine H. Graham, Anthony Lehmann, John Leathwick, and Simon Ferrier. Sample selection bias and presence-only distribution models: implications for background and pseudo-absence data. *Ecological Applications*, 19(1):181–197, January 2009. ISSN 1051-0761, 1939-5582. doi:10.1890/07-2153.1. URL <https://esajournals.onlinelibrary.wiley.com/doi/10.1890/07-2153.1>.
- Ian W Renner and David I Warton. Equivalence of maxent and poisson point process models for species distribution modeling in ecology. *Biometrics*, 69(1):274–281, 2013.
- Christophe Botella, Alexis Joly, Pascal Monestiez, Pierre Bonnet, and François Munoz. Bias in presence-only niche models related to sampling effort and species niches: Lessons for background point selection. *PLOS ONE*, 15(5):e0232078, May 2020. ISSN 1932-6203. doi:10.1371/journal.pone.0232078. URL <https://dx.plos.org/10.1371/journal.pone.0232078>.
- Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep learning*. MIT press, 2016.
- Kurt Hornik, Maxwell Stinchcombe, and Halbert White. Multilayer feedforward networks are universal approximators. *Neural networks*, 2(5):359–366, 1989. URL <https://www.sciencedirect.com/science/article/pii/0893608089900208>. Publisher: Elsevier.
- Yann LeCun, Bernhard Boser, John S. Denker, Donnie Henderson, Richard E. Howard, Wayne Hubbard, and Lawrence D. Jackel. Backpropagation applied to handwritten zip code recognition. *Neural computation*, 1(4):541–551, 1989. URL <https://ieeexplore.ieee.org/abstract/document/6795724/>. Publisher: MIT Press.
- Steven J Phillips and Miroslav Dudík. Modeling of species distributions with maxent: new extensions and a comprehensive evaluation. *Ecography*, 31(2):161–175, 2008.
- Osamu Komori, Yusuke Saigusa, Shinto Eguchi, and Yasuhiro Kubota. Cumulant-based approximation for fast and efficient prediction for species distribution, May 2024. URL <http://arxiv.org/abs/2405.14456>. arXiv:2405.14456.
- Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. Deep learning. *nature*, 521(7553):436–444, 2015. URL [https://idp.nature.com/authorize/casa?redirect\\_uri=https://www.nature.com/articles/nature14539&casa\\_token=dxn0tKfQj1sAAAAA:f3neASfT5epH41D2TIaCQf67m0jqF2yS17h-gNuj0iOBH\\_I0NbmDoWh0A\\_Ao0QrDTrpEnLnrPnJnjKjYJE](https://idp.nature.com/authorize/casa?redirect_uri=https://www.nature.com/articles/nature14539&casa_token=dxn0tKfQj1sAAAAA:f3neASfT5epH41D2TIaCQf67m0jqF2yS17h-gNuj0iOBH_I0NbmDoWh0A_Ao0QrDTrpEnLnrPnJnjKjYJE). Publisher: Nature Publishing Group UK London.



- Jürgen Schmidhuber. Deep learning in neural networks: An overview. *Neural networks*, 61:85–117, 2015. URL <https://www.sciencedirect.com/science/article/pii/S0893608014002135>. Publisher: Elsevier.
- Di Chen, Yexiang Xue, Daniel Fink, Shuo Chen, and Carla P Gomes. Deep multi-species embedding. In *Proceedings of the 26th International Joint Conference on Artificial Intelligence*, pages 3639–3646, 2017.
- Christophe Botella, Alexis Joly, Pierre Bonnet, Pascal Monestiez, and François Munoz. A deep learning approach to species distribution modelling. *Multimedia Tools and Applications for Environmental & Biodiversity Informatics*, pages 169–199, 2018.
- Benjamin Kellenberger, Kevin Winner, and Walter Jetz. The Performance and Potential of Deep Learning for Predicting Species Distributions, August 2024. URL <http://biorxiv.org/lookup/doi/10.1101/2024.08.09.607358>.
- Joaquim Estopinan, Pierre Bonnet, Maximilien Servajean, François Munoz, and Alexis Joly. Modelling Species Distributions with Deep Learning to Predict Plant Extinction Risk and Assess Climate Change Impacts, January 2024. URL <http://arxiv.org/abs/2401.05470>. arXiv:2401.05470 [cs, q-bio, stat].
- Robin Zbinden, Nina Van Tiel, Benjamin Kellenberger, Lloyd Hughes, and Devis Tuia. On the Selection and Effectiveness of Pseudo-Absences for Species Distribution Modeling with Deep Learning, 2024. URL <https://www.ssrn.com/abstract=4684222>.
- Ian W. Renner, Jane Elith, Adrian Baddeley, William Fithian, Trevor Hastie, Steven J. Phillips, Gordana Popovic, and David I. Warton. Point process models for presence-only analysis. *Methods in Ecology and Evolution*, 6(4): 366–379, April 2015. ISSN 2041-210X, 2041-210X. doi:10.1111/2041-210X.12352. URL <https://besjournals.onlinelibrary.wiley.com/doi/10.1111/2041-210X.12352>. Number: 4.
- Noel AC Cressie. Statistics for spatial data. John Willy and Sons. *Inc., New York*, 800, 1993.
- Diederik P Kingma. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- Prannay Khosla, Piotr Teterwak, Chen Wang, Aaron Sarna, Yonglong Tian, Phillip Isola, Aaron Maschiot, Ce Liu, and Dilip Krishnan. Supervised contrastive learning. *Advances in neural information processing systems*, 33: 18661–18673, 2020.
- Roosbeh Valavi, Jane Elith, José J. Lahoz-Monfort, and Gurutzeta Guillera-Arroita. blockCV: An r package for generating spatially or environmentally separated folds for k-fold cross-validation of species distribution models. *Methods in Ecology and Evolution*, 10(2):225–232, 2019. ISSN 2041-210X. doi:10.1111/2041-210X.13107. URL <https://onlinelibrary.wiley.com/doi/abs/10.1111/2041-210X.13107>. \_eprint: <https://onlinelibrary.wiley.com/doi/pdf/10.1111/2041-210X.13107>.
- David R. Roberts, Volker Bahn, Simone Ciuti, Mark S. Boyce, Jane Elith, Gurutzeta Guillera-Arroita, Severin Hauenstein, José J. Lahoz-Monfort, Boris Schröder, Wilfried Thuiller, David I. Warton, Brendan A. Wintle, Florian Hartig, and Carsten F. Dormann. Cross-validation strategies for data with temporal, spatial, hierarchical, or phylogenetic structure. *Ecography*, 40(8):913–929, 2017. ISSN 1600-0587. doi:10.1111/ecog.02881. URL <https://onlinelibrary.wiley.com/doi/abs/10.1111/ecog.02881>. \_eprint: <https://onlinelibrary.wiley.com/doi/pdf/10.1111/ecog.02881>.
- Philipp Brun, Dirk N Karger, Damaris Zurell, Patrice Descombes, Lucienne C de Witte, Riccardo de Lutio, Jan Dirk Wegner, and Niklaus E Zimmermann. Multispecies deep learning using citizen science data produces more informative plant community models. *Nature Communications*, 15(1):4421, 2024.
- Donald J Benkendorf and Charles P Hawkins. Effects of sample size and network depth on a deep learning approach to species distribution modeling. *Ecological Informatics*, 60:101137, 2020.
- Tyler E. Schartel and Yong Cao. Background selection complexity influences Maxent predictive performance in freshwater systems. *Ecological Modelling*, 488:110592, February 2024. ISSN 0304-3800. doi:10.1016/j.ecolmodel.2023.110592. URL <https://www.sciencedirect.com/science/article/pii/S0304380023003228>.
- Jimmy Ba and Rich Caruana. Do deep nets really need to be deep? *Advances in neural information processing systems*, 27, 2014. URL <https://proceedings.neurips.cc/paper/2014/hash/ea8fcd92d59581717e06eb187f10666d-Abstract.html>.
- Maithra Raghu, Ben Poole, Jon Kleinberg, Surya Ganguli, and Jascha Sohl-Dickstein. On the expressive power of deep neural networks. In *international conference on machine learning*, pages 2847–2854. PMLR, 2017. URL <https://proceedings.mlr.press/v70/raghu17a.html>.
- Cory Merow, Mathew J. Smith, Thomas C. Edwards Jr, Antoine Guisan, Sean M. McMahon, Signe Normand, Wilfried Thuiller, Rafael O. Wüest, Niklaus E. Zimmermann, and Jane Elith. What do we gain from simplicity versus complexity in species distribution models? *Ecography*, 37(12):1267–1281, 2014. ISSN 1600-

0587. doi:10.1111/ecog.00845. URL <https://onlinelibrary.wiley.com/doi/abs/10.1111/ecog.00845>.  
\_eprint: <https://onlinelibrary.wiley.com/doi/pdf/10.1111/ecog.00845>.

Vira Koshkina, Yan Wang, Ascelin Gordon, Robert M Dorazio, Matt White, and Lewi Stone. Integrated species distribution models: combining presence-background data and site-occupancy data with imperfect detection. *Methods in Ecology and Evolution*, 8(4):420–430, 2017.

Robert M Dorazio. Accounting for imperfect detection and survey bias in statistical analysis of presence-only data. *Global Ecology and Biogeography*, 23(12):1472–1484, 2014.

David AW Miller, Krishna Pacifici, Jamie S Sanderlin, and Brian J Reich. The recent past and promising future for data integration methods to estimate species' distributions. *Methods in Ecology and Evolution*, 10(1):22–37, 2019.

Nick JB Isaac, Marta A Jarzyna, Petr Keil, Lea I Dambly, Philipp H Boersch-Supan, Ella Browning, Stephen N Freeman, Nick Golding, Gurutzeta Guillera-Arroita, Peter A Henrys, et al. Data integration for large-scale models of species distributions. *Trends in ecology & evolution*, 35(1):56–67, 2020.

Philip S Mostert and Robert B O'Hara. Pointedsdms: An r package to help facilitate the construction of integrated species distribution models. *Methods in Ecology and Evolution*, 14(5):1200–1207, 2023.

## A Supplementary theoretical aspects

### A.1 Relationship between Poisson loss and Maxent loss

For this normalisation,  $\lambda$  and  $y$  at site  $i$  for a given species  $j$  are expressed as the normalised values  $\tilde{\lambda}_{ij}$  and  $\tilde{y}_{ij}$

$$\tilde{\lambda}_{ij} = \frac{\lambda_{ij}}{\sum_{k=1}^K \lambda_{kj}} \quad (11)$$

$$\tilde{y}_{ij} = \frac{y_{ij}}{\sum_{k=1}^K y_{kj}} \quad (12)$$

Thus the Poisson loss in a context where intensities and the number of occurrences are normalised per site can be formulated as follow:

$$\mathcal{L}(\tilde{\lambda}, \tilde{y}) = \frac{1}{KN} \sum_{i=1}^K \sum_{j=1}^N \left( \tilde{\lambda}_{ij} - \tilde{y}_{ij} \log \tilde{\lambda}_{ij} \right) \quad (13)$$

$$= \frac{1}{KN} \sum_{i=1}^K \sum_{j=1}^N \left( \frac{\lambda_{ij}}{\sum_{k=1}^K \lambda_{kj}} - \frac{y_{ij}}{\sum_{k=1}^K y_{kj}} \log \left( \frac{\lambda_{ij}}{\sum_{k=1}^K \lambda_{kj}} \right) \right) \quad (14)$$

$$= \frac{1}{KN} \sum_{i=1}^K \sum_{j=1}^N \left( \frac{\lambda_{ij}}{\sum_{k=1}^K \lambda_{kj}} \right) - \frac{1}{KN} \sum_{i=1}^K \sum_{j=1}^N y_{ij} \log \left( \frac{\lambda_{ij}}{\sum_{k=1}^K \lambda_{kj}} \right) \quad (15)$$

$$= \frac{1}{K} - \frac{1}{KN} \sum_{i=1}^K \sum_{j=1}^N \left( \frac{y_{ij}}{\sum_{k=1}^K y_{kj}} \right) \log \left( \frac{\lambda_{ij}}{\sum_{k=1}^K \lambda_{kj}} \right) \quad (16)$$

$$= \frac{1}{K} + \mathcal{L}_{\mathcal{H}}(\lambda, y) \quad (17)$$

### A.2 The case of Maxent

Maximum entropy (Maxent) is a widely used method in species distribution modelling (SDM) that estimates the distribution of a species on the basis of environmental variables. The objective of Maxent is to approximate the probability distribution  $p(x)$  of species occurrences over a set of sites using PO data. The method maximises the entropy of the predicted distribution, subject to environmental constraints derived from occurrence data.

The Maxent approach can be seen as a special case of a Poisson Point Process (PPP). A PPP models the occurrence of species as events in space, with a function  $\lambda(x)$  describing the intensity of occurrences at location  $x$ . If we denote by  $X$  the set of all possible sites, the probability density of species occurrence can be written as follows:

$$p(x) = \frac{\lambda(x)}{\int_X \lambda(x') dx'} \quad (18)$$

Maxent is equivalent to estimating the intensity  $\lambda(x)$  using the environmental variables at each site  $x$ . In Maxent,  $\lambda(x)$  is modelled as a log-linear function of features (e.g. environmental covariates):

$$\lambda(x) = \exp \left( \beta_0 + \sum_{i=1}^p \beta_i f_i(x) \right) \quad (19)$$

where  $f_i(x)$  are the environmental characteristics at location  $x$ , and  $\beta_i$  are the parameters to be learned. The objective is to estimate  $\beta_i$  by maximising the probability of observed occurrences. Given a set of occurrence points  $S = \{x_1, x_2, \dots, x_n\}$ , the loss function of Maxent is to maximise the likelihood under the assumption of a Poisson process. The complete Maxent loss can be written as follows:

$$\mathcal{L}(\beta) = -\frac{1}{n} \sum_{x \in S} \log p(x) + \lambda \sum_{i=1}^p |\beta_i| \quad (20)$$

where  $p(x)$  is the probability of occurrence at site  $x$ ,  $\lambda$  is a regularisation parameter, and  $|\beta_i|$  is the L1 regularisation term that promotes sparsity in the learned coefficients. In Maxent, the partition function  $Z(\beta)$  is used to normalise the predicted distribution so that the sum of the probabilities is equal to 1. It is defined by the following integral:

$$Z(\beta) = \int_X \exp \left( \beta_0 + \sum_{i=1}^p \beta_i f_i(x) \right) dx \quad (21)$$

The predicted probability of occurrence  $p(x)$  is then:

$$p(x) = \frac{\exp(\beta_0 + \sum_{i=1}^p \beta_i f_i(x))}{Z(\beta)} \quad (22)$$

### A.3 Minimising the loss on each batch implies minimising the full loss

In DeepMaxent, the term batch is not to be understood in its classical sense in machine learning or statistics. This is because the formula of the loss for a batch is not simply the restriction of the full loss to the terms of that batch, due to the partition functions. Thus, unlike in a classical setting, using the SGD with this type of batches might not necessarily lead to an estimator that is good regarding the full loss. To address this, we show below that an estimator that minimises all the batch-wise losses actually minimises the full loss, thus justifying the use of the SGD.

**Notations:** Note for any  $n \in \mathbb{N}^*$ ,  $\Delta^n = \{t_0, \dots, t_n \in \mathbb{R}^{n+1} \mid \sum_{i=0}^n t_i = 1, \forall i, t_i \geq 0\}$  the  $n$ -probability simplex, i.e. the space of probability distributions on a set of  $n + 1$  elements. Without loss of generality, we consider a single species with occurrence (pseudo-)count  $y_k$  in site  $k$ .  $\forall k \in [1, K], y_k > 0$ , consistently with our implementation, where  $0 < \delta \ll 1$  is added to each raw occurrence count to avoid numerical problems. For any  $B \subset [1, K]$ , we note  $y_B := \{y_k\}_{k \in B}$ , and also apply this notation to  $\lambda_B$ . We express the batch cross-entropy loss function of DeepMaxent with our notation in equation 23.

$$\mathcal{L}(\lambda_B, y_B) = -\frac{1}{|B|} \sum_{k \in B} \frac{y_k}{\sum_{i \in B} y_i} \log \left( \frac{\lambda_k}{\sum_{i \in B} \lambda_i} \right) \quad (23)$$

According to Gibbs inequality, we have that  $\operatorname{argmin}_{p \in \Delta^{|B|-1}} \mathcal{L}(p, y_B) = \{y_k / \sum_{i \in B} y_i\}_{k \in B}$ .

Now, assume that an estimator of the intensity  $\lambda_{[1, K]}^*$  ( $\forall k \in [1, K], \lambda_k^* > 0$ ) minimises our loss for any batch of size  $n$  ( $1 < n < K$ ). Formally,  $\lambda^*$  fulfils the property  $P(n)$ :  $\forall B \subset [1, K]$  such that  $|B| = n$ ,  $\{\lambda_k^* / \sum_{i \in B} \lambda_i^*\}_{k \in B} = \{y_k / \sum_{i \in B} y_i\}_{k \in B}$

As a preliminary corollary, we have that  $P(n) \Rightarrow P'(n)$ :  $\forall B_1, B_2 \subset [1, K] / |B_1| = |B_2| = n$  and  $B_1 \cap B_2 \neq \emptyset$ ,  $\frac{\sum_{i \in B_1} \lambda_i^*}{\sum_{i \in B_2} \lambda_i^*} = \frac{\sum_{i \in B_1} y_i}{\sum_{i \in B_2} y_i}$ , because from  $P(n)$  we get  $\frac{\lambda_k^*}{y_k} = \frac{\sum_{i \in B_1} \lambda_i^*}{\sum_{i \in B_1} y_i} = \frac{\sum_{i \in B_2} \lambda_i^*}{\sum_{i \in B_2} y_i}$  which leads to the result.

**Property:**  $P(n) \Rightarrow P(K)$ , that is, if  $\lambda^*$  fulfils  $P(n)$ , it minimises as well the full loss  $\mathcal{L}(\lambda_{[1, K]}^*, y_{[1, K]})$ .

**Proof:** Let's show that  $P(n) \Rightarrow P(n + 1)$ , implying  $P(K)$  by induction.

$P(n + 1) \Leftrightarrow \forall B \subset [1, K], l \in [1, K], k \notin B$ , we have  $\{\lambda_k^* / \sum_{i \in B \cup \{l\}} \lambda_i^*\}_{k \in B \cup \{l\}} = \{y_k / \sum_{i \in B \cup \{l\}} y_i\}_{k \in B \cup \{l\}}$

Let be such  $B$  and  $l$ , then  $\forall i, j \in B, i \neq j$ ,

$$\begin{aligned}
\frac{\lambda_i^*}{\sum_{k \in B \cup I} \lambda_k^*} &= \frac{\lambda_i^*}{\sum_{k \in B \cup I \setminus j} \lambda_k^* + \lambda_j^*} \\
&= \frac{\lambda_i^* / \sum_{k \in B} \lambda_k^*}{\frac{\sum_{k \in B \cup I \setminus j} \lambda_k^*}{\sum_{k \in B} \lambda_k^*} + \frac{\lambda_j^*}{\sum_{k \in B} \lambda_k^*}} \\
&= \frac{y_i / \sum_{k \in B} y_k}{P(n), P'(n) \frac{\sum_{k \in B \cup I \setminus j} y_k}{\sum_{k \in B} y_k} + \frac{y_j}{\sum_{k \in B} y_k}} \\
&= \frac{y_i}{\sum_{k \in B \cup I \setminus j} y_k + y_j} \\
&= \frac{y_i}{\sum_{k \in B \cup I} y_k}
\end{aligned} \tag{24}$$

So  $\lambda^*$  satisfies  $P(n+1)$ .

## B Dataset description

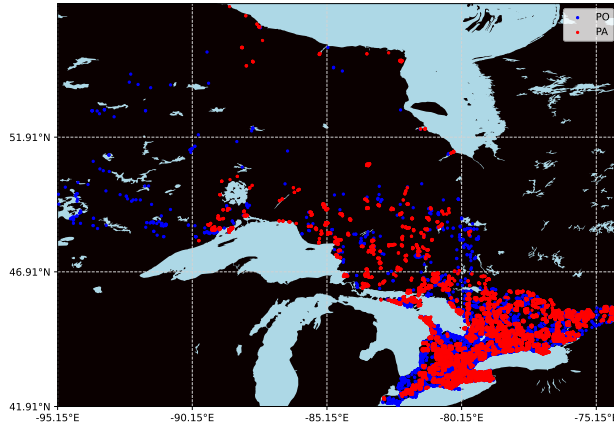


Figure 3: Occurrences for all 20 bird species in Canada, for PO and PA data)

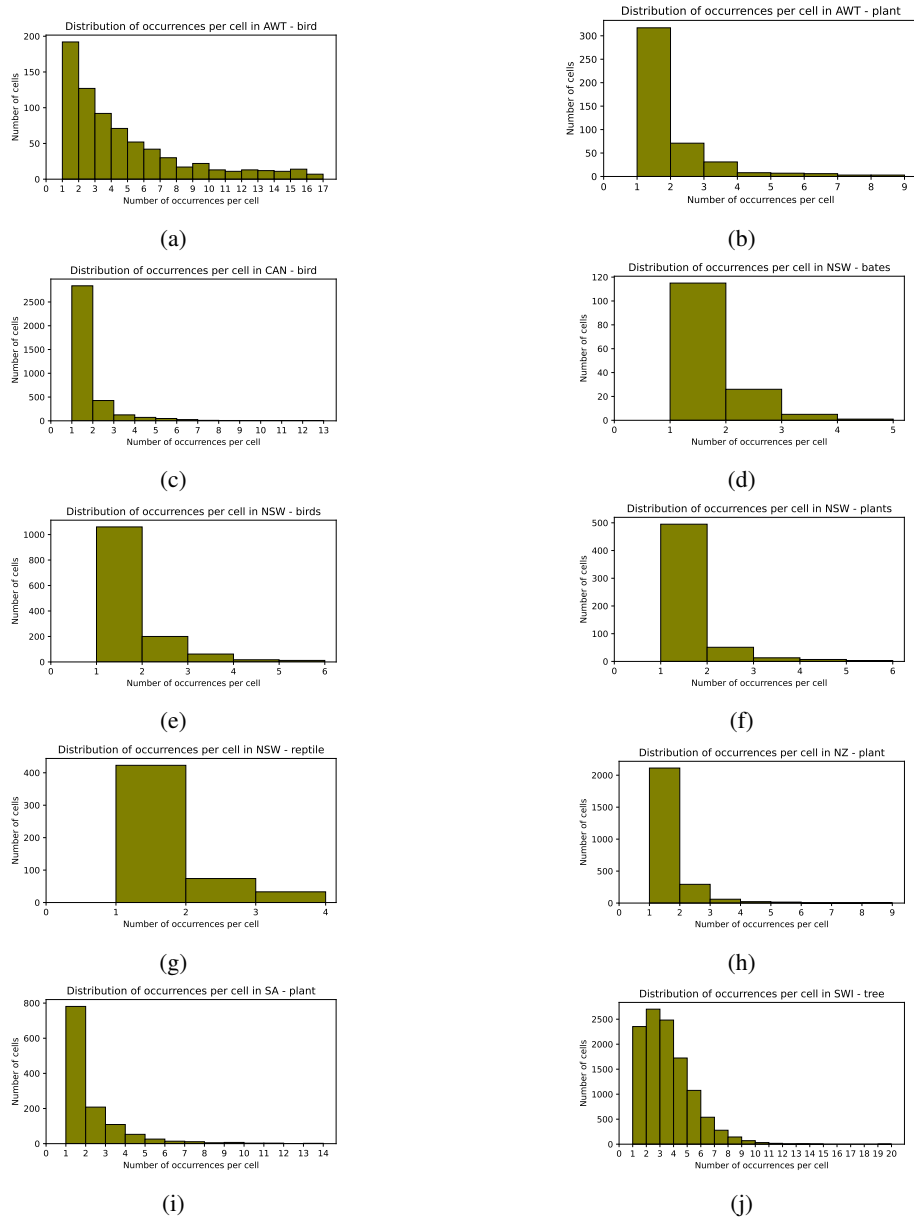


Figure 4: Distribution of occurrence numbers in cells that contain at least one occurrence for each region and biological group: (a) AWT bird, (b) AWT plant, (c) CAN bird, (d) NSW bates, (e) NSW bird, (f) NSW plant, (g) NSW reptile, (h) NZ plant, (i) SA plant and (j) SWI tree

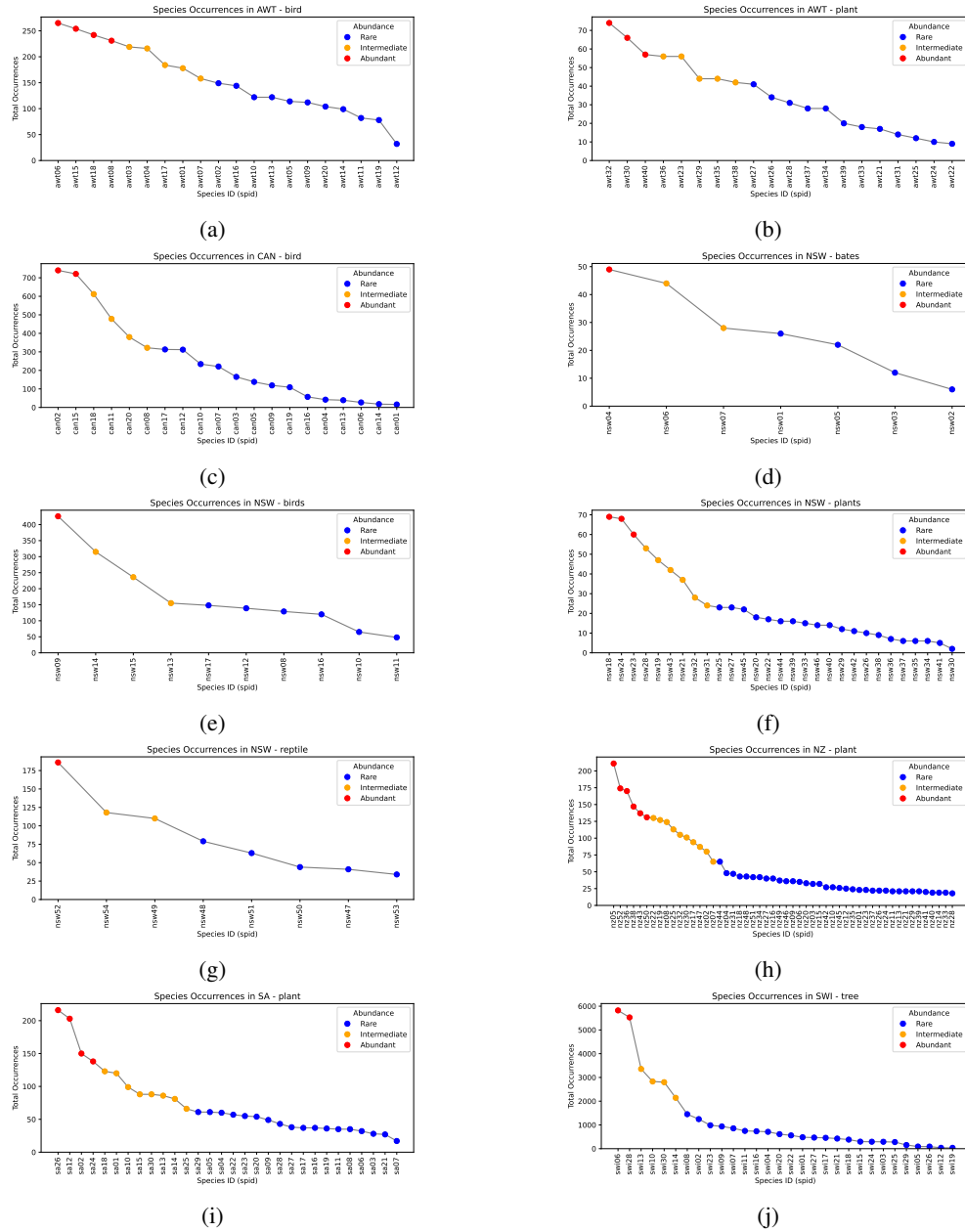


Figure 5: Species occurrence numbers in PO data for each region and biological group: (a) AWT bird, (b) AWT plant, (c) CAN bird, (d) NSW bates, (e) NSW bird, (f) NSW plant, (g) NSW reptile, (h) NZ plant, (i) SA plant and (j) SWI tree. It's displayed in descending order. Colours correspond to abundance classes: Common, Intermediate, and Rare

## C A suggestion for cross-validating DeepMaxent using PO data

In this study, model calibration was performed using presence-only (PO) data, whereas model evaluation used presence-absence (PA) data with a specific metric, the Area Under the Curve (AUC). This difference introduces a fundamental challenge in choosing the appropriate evaluation metric in the cross-validation step using PO data. PO and PA data are collected under different conditions with bias related to inhomogeneous sampling effort in the case of PO data. To fine-tune the DeepMaxent model’s hyperparameters, spatial blocking was employed for cross-validation, drawing on geographical data to optimise the model’s performance. The cross-validation was performed following the approach of Zbinden et al. [2024], Roberts et al. [2017]. In more detail, we performed a spatially stratified 10-fold cross-validation using a grid of 25 spatial blocks. Each fold used a unique subset of blocks for validation, distinct from previous folds, with blocks selected randomly yet balanced to ensure an even distribution of presence data across folds.

Table 4: Parameters and Tested Values for cross-validation.

Parameter	Tested Values
Batch Size	250, 2500
Hidden Layers	1, 2, 3
Weight Decay	$10^{-3}$ , $10^{-2}$ , $10^{-1}$

To achieve this, all combinations of hyperparameter values visible in Table 4 were evaluated. The evaluation criterion used was the average AUC across all validation sets. The optimal model was obtained with two hidden layers, a learning rate of 0.0002, a weight decay of  $10^{-2}$  and a batch size of 250.

Table 5: All combinations of hyperparameter values and their cross-validation results on PO data.

Batch Size	Hidden Layers	Weight Decay	AUC_CV
250	1	0.0001	0.632
250	1	0.001	0.631
250	1	0.01	0.649
250	2	0.0001	0.654
250	2	0.001	0.659
250	2	0.01	<b>0.673</b>
250	3	0.0001	0.619
250	3	0.001	0.624
250	3	0.01	0.635
2500	1	0.0001	0.634
2500	1	0.001	0.631
2500	1	0.01	0.632
2500	2	0.0001	0.659
2500	2	0.001	0.659
2500	2	0.01	0.660
2500	3	0.0001	0.617
2500	3	0.001	0.619
2500	3	0.01	0.634

To obtain optimal hyperparameter values, one of the main challenges is to find the right combination of hyperparameters during cross-validation while maintaining computational efficiency. Techniques such as subsampling the dataset or limiting the training process to a reduced number of epochs can help to reduce the computational time for each configuration.

## D Supplementary results

### D.1 Batch size effect on model performance

The impact of batch size on the average AUC values obtained by the DeepMaxent method in the different regions of the dataset is illustrated in Table 6.



Table 6: Impact of batch size on DeepMaxent performance across different regions. The table presents the accuracy scores for each region (AWT, CAN, NSW, NZ, SA, SWI) and the average (avg) performance across all regions.

Batch size	regions						avg
	AWT	CAN	NSW	NZ	SA	SWI	
10	0.702	<u>0.732</u>	0.747	<b>0.755</b>	<b>0.805</b>	0.848	0.765
25	0.705	<u>0.731</u>	0.748	<b>0.755</b>	<b>0.805</b>	0.849	0.765
100	0.711	0.730	<b>0.754</b>	<b>0.755</b>	<u>0.804</u>	<u>0.849</u>	<b>0.767</b>
250	0.714	<u>0.732</u>	<u>0.752</u>	<u>0.754</u>	0.803	<b>0.850</b>	<b>0.767</b>
1000	<u>0.715</u>	<b>0.733</b>	0.750	0.749	0.801	0.848	<u>0.766</u>
2500	<b>0.716</b>	<b>0.733</b>	0.748	0.742	0.800	0.847	0.764

The mean AUC values for small batch sizes such as batch size 10 averaged 0.765. Increasing the batch size to 25 does not improve the overall AUC. Even though very small, the improvement comes mainly from AWT, NSW, and SWI. The average AUC in all regions shows better values for batch sizes of 100 and 250, with values of 0.767. Conversely, when the batch size increases to 1000 or 2500, the average AUC values for all regions decrease, reaching 0.766 and 0.764, respectively.

When analysing specific regions, we notice that regions like AWT and NSW show high variability in their AUC values. AWT exhibits the most instability, as its AUC value fluctuates between 0.702 and 0.716 across different batch sizes, suggesting that it may require a larger batch size for consistent performance. Conversely, SWI maintains a high AUC value around 0.848 with minimal changes, indicating that it is relatively unaffected by batch size variations. Similarly, regions such as CAN, NSW and SA also show limited impact, with AUC values remaining close to 0.732, 0.751 and 0.801, respectively, regardless of the batch size.

For some regions, such as SWI and CAN, the effects of batch size appear to be minimal, with good approximations obtained regardless of batch size. However, a medium batch size can improve computational efficiency on machines, providing a balance between performance and resource utilisation compared to larger batch sizes. On the other hand, the AWT remains the most unstable region, indicating that it requires more precise parametrisation to achieve optimal performance.

## D.2 Influence of hidden layer number on model performance

Table 7 presents the impact of hidden layer number on the performance of the DeepMaxent model, measured by AUC values across regions: AWT, CAN, NSW, NZ, SA, and SWI.

Table 7: Impact of hidden layer number on DeepMaxent performance with L2 regularisation ( $w=1e-2$ ) across different regions. The table presents AUC value for each region (AWT, CAN, NSW, NZ, SA, SWI) and the average (avg) AUC across all regions.

Hidden layer number	regions						avg
	AWT	CAN	NSW	NZ	SA	SWI	
1	<b>0.719</b>	<b>0.735</b>	<b>0.753</b>	0.751	0.800	0.847	<b>0.767</b>
2	<u>0.714</u>	<u>0.732</u>	<u>0.752</u>	<b>0.754</b>	0.803	<b>0.850</b>	<b>0.767</b>
3	<u>0.710</u>	<u>0.729</u>	<u>0.752</u>	<u>0.752</u>	<u>0.805</u>	<b>0.850</b>	<u>0.766</u>
4	0.705	0.727	0.748	0.752	<u>0.805</u>	<u>0.849</u>	0.764
5	0.703	0.722	0.743	0.749	<b>0.806</b>	0.848	0.762
6	0.698	0.719	0.739	0.747	<b>0.806</b>	0.847	0.759

In some regions, the AUC values decrease as the number of hidden layers increases. For instance, in AWT, CAN and NSW, the highest AUC values are achieved with just one hidden layer. In contrast, for regions like SA, increasing the number of layers leads to a slight improvement in AUC performance, reaching an AUC value of 0.806. In other regions, an optimum is reached, such as in SWI with an AUC of 0.850 and NZ with 0.754, both with two hidden layers. With one layer, the model achieves an average AUC of 0.767, the highest overall, particularly performing well in regions like AWT with a AUC of 0.720 and NSW with an AUC of 0.754.

The impact of the number of hidden layers varies by region. In some cases, the changes are minimal, such as for SWI, where the AUC fluctuates by only 0.003, while for the AWT dataset, the variations are more pronounced. Overall, using

one or two hidden layers results in the highest average AUC across all regions, at 0.767. However, with two hidden layers, the model produces the highest or second-highest scores in most regions, suggesting that this configuration may represent the optimal balance across all regions rather than favouring a specific one. This is further supported by the Pearson correlation values provided in the appendix (see Table 11).

### D.3 Sensitivity to the L2 regularisation weight

The average AUC values per region and across all regions for various weight decay values are shown in Table 8.

Table 8: Comparison of DeepMaxent performance according to different weight decay values for L2 regularisation. The criteria are the average AUC per region and the average AUC for all regions. The best average AUC for each column is highlighted in bold, while the second-best averaged AUC is underlined.

weight decay value	regions						
	AWT	CAN	NSW	NZ	SA	SWI	avg
0	0.713	0.719	0.750	0.744	0.804	0.843	0.762
1e-6	0.713	0.719	0.750	0.744	<b>0.804</b>	0.843	0.762
3e-6	0.713	0.719	0.750	0.744	<b>0.804</b>	0.843	0.762
1e-5	0.713	0.719	0.750	0.744	<b>0.804</b>	0.844	0.762
3e-5	0.713	0.721	0.750	0.746	<b>0.804</b>	0.845	0.763
1e-4	0.713	0.726	0.751	<u>0.750</u>	<b>0.804</b>	<u>0.848</u>	<u>0.765</u>
3e-4	0.714	<u>0.732</u>	<u>0.752</u>	<b>0.754</b>	<u>0.803</u>	<b>0.850</b>	<b>0.767</b>
1e-3	<u>0.718</u>	<b>0.734</b>	<b>0.753</b>	0.745	0.799	0.843	<u>0.765</u>
3e-3	<b>0.722</b>	<u>0.732</u>	0.749	0.723	0.780	0.833	0.757
1e-2	0.716	0.730	0.740	0.707	0.754	0.805	0.742
3e-2	0.691	0.705	0.715	0.662	0.737	0.780	0.715
1e-1	0.642	0.661	0.670	0.585	0.669	0.718	0.657
3e-1	0.592	0.636	0.615	0.560	0.604	0.644	0.609

Without regularisation (with  $w=0$ ), DeepMaxent method achieved an average AUC value of 0.762, highlighted already a strong performance in comparison of classic method (see Table 2). When L2 regularisation with a weight decay value of  $w = 1 \times 10^{-6}$  or  $w = 1 \times 10^{-5}$  is applied, the AUC values for each region, as well as the overall average AUC, remain approximatively unchanged. Marginal improvements on the order of 0.001 are observed for Switzerland. These results are therefore approximately the same as those without regularisation. This suggests that the weight decay value is too small to significantly affect the outcomes, resulting in minimal L2 regularisation.

When the weight decay value is increased to  $w = 3 \times 10^{-5}$ , slight changes are observed, with noticeable improvements particularly in CAN and SWI regions. The overall average AUC begins to be affected by regularisation, with an overall AUC value of 0.763. With a weight decay value of  $1 \times 10^{-4}$ , regularisation has a more pronounced effect on the results, yielding an overall average AUC of 0.765. Significant improvements are observed, particularly for CAN and SWI, with AUC values of 0.726 and 0.848, respectively. When weight decay has a value of  $3 \times 10^{-4}$ , the overall AUC is the highest with a value of 0.767. In addition, five of the six regions obtained the best or the second-best values in this ablation study. This is the case for CAN, NSW, NZ, SA and SWI with AUC values of 0.732, 0.752, 0.754, 0.803 and 0.850, respectively.

On the other hand, with higher weight decay values, the overall average AUC across all regions decreases. For a weight decay value of  $3 \times 10^{-2}$ , the mean AUC decreases slightly to reach a value of 0.742. All AUC values for all regions decrease. Finally, with higher weight decay values of  $1 \times 10^{-1}$  and  $3 \times 10^{-2}$ , the overall mean AUC decreases significantly.

In a multi-species Maxent model, applying a normalisation based on the sum of  $y$  (presences weighted by probability) for each species can indeed impact the distribution of relative presence probabilities between abundant and rare species. When a uniform weight is assigned to each species (for example, setting each species' weight to 1), there is a risk that abundant species will have their occurrences assigned lower probabilities, while, conversely, the occurrences of rare species may be prioritised.

Table 9: AUC values by bird species in Canada (Ontario) dataset according to different loss functions. The best average AUC for each species (row) is highlighted in bold, while the second-best AUC value is underlined.

spid	count	class	Batch size					
			10	25	100	250	1000	2500
can01	16	rare	0.943	0.942	0.941	0.942	0.942	0.942
can02	740	abundant	0.709	0.708	0.707	0.707	0.705	0.706
can03	165	rare	0.658	0.655	0.655	0.659	0.662	0.670
can04	42	rare	0.755	0.754	0.741	0.750	0.767	0.786
can05	138	rare	0.656	0.657	0.653	0.651	0.647	0.638
can06	27	rare	0.886	0.886	0.886	0.886	0.886	0.885
can07	221	rare	0.695	0.692	0.690	0.687	0.673	0.663
can08	322	common	0.552	0.554	0.556	0.557	0.555	0.547
can09	119	rare	0.731	0.733	0.735	0.735	0.724	0.714
can10	234	rare	0.700	0.700	0.703	0.703	0.696	0.691
can11	478	common	0.594	0.591	0.590	0.592	0.590	0.590
can12	312	rare	0.848	0.847	0.846	0.846	0.846	0.846
can13	39	rare	0.701	0.695	0.687	0.682	0.681	0.689
can14	18	rare	0.923	0.921	0.918	0.919	0.920	0.921
can15	721	abundant	0.622	0.618	0.627	0.636	0.638	0.644
can16	57	rare	0.858	0.853	0.844	0.843	0.848	0.853
can17	313	rare	0.786	0.789	0.789	0.794	0.799	0.801
can18	612	common	0.762	0.760	0.764	0.765	0.779	0.792
can19	109	rare	0.482	0.488	0.494	0.511	0.519	0.498
can20	380	common	0.778	0.779	0.777	0.776	0.778	0.782

Table 10: AUC values by bird species in Canada (Ontario) dataset according to different weight decay values.

spid	count	class	w=0	w=1e-6	w=3e-6	w=1e-5	w=3e-5	w=1e-4	w=3e-4	w=1e-3	w=3e-3	w=1e-2	w=3e-2	w=1e-1	w=3e-1
can01	16	rare	0.941	0.941	0.941	0.941	0.942	0.942	0.943	0.942	0.943	0.943	0.895	0.729	0.541
can02	740	abundant	0.706	0.706	0.705	0.705	0.705	0.706	0.706	0.708	0.716	0.721	0.717	0.702	0.688
can03	165	rare	0.652	0.652	0.653	0.653	0.654	0.656	0.661	0.673	0.674	0.667	0.662	0.627	0.579
can04	42	rare	0.728	0.728	0.731	0.731	0.734	0.743	0.764	0.801	0.783	0.765	0.710	0.547	0.544
can05	138	rare	0.649	0.650	0.649	0.649	0.650	0.649	0.648	0.639	0.625	0.621	0.611	0.554	0.520
can06	27	rare	0.887	0.887	0.887	0.887	0.886	0.886	0.886	0.886	0.884	0.883	0.845	0.736	0.666
can07	221	rare	0.683	0.683	0.682	0.682	0.682	0.682	0.678	0.663	0.657	0.662	0.662	0.621	0.602
can08	322	common	0.562	0.562	0.562	0.562	0.561	0.560	0.555	0.539	0.525	0.521	0.526	0.526	0.522
can09	119	rare	0.739	0.739	0.738	0.738	0.737	0.735	0.728	0.704	0.689	0.699	0.712	0.674	0.604
can10	234	rare	0.705	0.705	0.705	0.704	0.704	0.702	0.697	0.687	0.679	0.668	0.656	0.611	0.556
can11	478	common	0.590	0.590	0.590	0.590	0.591	0.591	0.591	0.592	0.595	0.595	0.592	0.575	0.558
can12	312	rare	0.844	0.844	0.844	0.844	0.845	0.846	0.846	0.848	0.850	0.848	0.844	0.829	0.775
can13	39	rare	0.651	0.652	0.652	0.652	0.656	0.665	0.682	0.712	0.734	0.732	0.690	0.621	0.581
can14	18	rare	0.902	0.902	0.903	0.904	0.907	0.913	0.919	0.922	0.923	0.921	0.772	0.582	0.555
can15	721	abundant	0.618	0.617	0.620	0.620	0.624	0.629	0.640	0.656	0.675	0.696	0.674	0.605	0.568
can16	57	rare	0.825	0.825	0.826	0.826	0.830	0.836	0.847	0.856	0.856	0.854	0.817	0.759	0.611
can17	313	rare	0.779	0.780	0.781	0.781	0.785	0.791	0.799	0.800	0.788	0.770	0.757	0.708	0.630
can18	612	common	0.760	0.760	0.761	0.761	0.762	0.764	0.775	0.795	0.799	0.798	0.789	0.759	0.690
can19	109	rare	0.534	0.535	0.535	0.536	0.536	0.533	0.521	0.474	0.461	0.469	0.481	0.483	0.487
can20	380	common	0.776	0.776	0.776	0.776	0.776	0.775	0.776	0.784	0.780	0.768	0.758	0.709	0.621

Table 11: Impact of hidden layer number on DeepMaxent performance with L2 regularisation ( $w=1e-2$ ). The table presents average AUC and Pearson coefficient values across all regions.

Hidden layer number	avg AUC	avg Pearson Coefficient
1	<b>0.767</b>	0.244
2	<b>0.767</b>	<b>0.247</b>
3	<b>0.767</b>	0.242
4	<u>0.764</u>	0.234
5	0.762	0.228
6	0.759	0.222

Table 12: Impact of batch size on DeepMaxent performance across all regions. The table presents average AUC and Pearson coefficient values across all regions.

Batch size	avg AUC	avg Pearson Coefficient
10	0.765	0.231
25	0.765	0.236
100	<b>0.767</b>	0.244
250	<b>0.767</b>	<b>0.247</b>
1000	<u>0.766</u>	<u>0.245</u>
2500	0.764	0.241

Table 13: Comparison of method performance of DeepMaxent using TGB and using a equal weight between species by region-averaged AUC and averaged over all regions. The best average AUC for each column is highlighted in bold, while the second-best averaged AUC is underlined.

weight decay value	regions						avg
	AWT	CAN	NSW	NZ	SA	SWI	
0	0.717	0.728	0.737	0.745	0.804	0.844	0.762
1e-6	0.716	0.725	0.742	0.745	0.804	0.844	0.763
3e-6	0.716	0.726	0.742	0.746	0.804	0.844	0.763
1e-5	0.716	0.727	0.742	0.746	0.805	0.846	0.763
3e-5	0.717	0.729	0.741	0.746	0.804	0.845	0.764
1e-4	0.716	0.731	0.740	0.745	0.799	0.839	0.761
3e-4	0.711	0.730	0.735	0.738	0.782	0.826	0.754
1e-3	0.687	0.727	0.716	0.707	0.755	0.795	0.731
3e-3	0.629	0.700	0.681	0.682	0.708	0.776	0.696

Table 14: Comparison of performance using DeepMaxent without TGB approach according to different weight decay values. The highest average AUC for each column is highlighted in bold, with the second-highest average AUC underlined.

Weight decay value	regions						avg
	AWT	CAN	NSW	NZ	SA	SWI	
1e-1	0.691	0.455	0.594	0.679	0.746	0.647	0.635
1e-2	0.665	0.469	0.691	0.719	0.773	0.770	0.681
1e-3	0.654	0.593	0.718	0.744	0.803	0.810	0.720
1e-4	0.652	0.601	0.713	0.731	0.802	0.803	0.717
1e-5	0.651	0.596	0.712	0.729	0.802	0.799	0.715
1e-6	0.652	0.595	0.712	0.730	0.802	0.798	0.715