

Toward a Monte Carlo approach to selecting climate variables in MaxEnt:

A case study using Cassin's Sparrow (*Peucaea cassinii*)

John L. Schnase*, Mark L. Carroll, Roger L. Gill, Glenn S. Tamkin,
Jian Li, Savannah L. Strong, Thomas P. Maxwell, and Mary E. Aronne

Office of Computational and Information Sciences and Technology,
NASA Goddard Space Flight Center, Greenbelt, Maryland, 20708, USA

*Corresponding author

Email: john.l.schnase@nasa.gov

Author contributions

JLS: Conceptualization, Formal analysis, Investigation, Methodology,
Writing – original draft. MLC: Conceptualization, Formal analysis, Methodology,
Writing – review & editing. RLG: Software. GST: Software. JL: Software.
SLS: Validation. TPM: Software. MEA: Visualization.

Abstract

MaxEnt is an important aid in understanding the influence of climate change on species distributions and abundance. There is growing interest in using IPCC-class global climate model outputs as environmental predictors in this work. These models provide realistic, global representations of the climate system, projections for hundreds of variables (including Essential Climate Variables), and combine observations from an array of satellite, airborne, and *in-situ* sensors. Unfortunately, direct use of this important class of data in MaxEnt modeling has been limited due to the large size of climate model output collections. In this study, we investigated the potential of a Monte Carlo method to find a useful subset of predictors in a larger collection of environmental variables in a reasonable amount of time. Our proposed solution takes an ensemble approach wherein many MaxEnt runs, each drawing on a small random subset of variables, converges on a global estimate of the top contributing subset of variables in the larger collection. The Monte Carlo approach resulted in a consistent set of top six variables within 540 runs, and the four most contributory variables of the top six accounted for approximately 93% of overall permutation importance in the final model. These preliminary results suggest that a Monte Carlo approach could offer a viable means of selecting environmental predictors for MaxEnt models that is amenable to parallelization and scalable to large data sets, including externally-stored collections. This points to the possibility of near-real-time multiprocessor implementations that could enable broader and more exploratory use of global climate model outputs in environmental niche modeling and aid in the discovery of viable predictors.

Introduction

MaxEnt is an important aid in understanding the influence of climate change on species distributions and abundance. Based on a machine learning approach to maximum entropy modeling, the software allows researchers to construct ecological niche models (ENMs) that estimate the habitat suitability of a species using occurrence data and a set of environmental variables [1–3]. The need for reliable climate projections in this work is leading to greater use of global climate model (GCM) outputs as predictors [4]. While creating important new opportunities for research, this trend is also creating a “Big Data” challenge for the MaxEnt community [5]. The largest and most sophisticated GCMs — sometimes referred to as “IPCC-class” models because of the critical role they play in the work of the Intergovernmental Panel on Climate Change (IPCC) — produce petabyte-scale data sets comprising hundreds of variables, a volume that vastly exceeds what is generally used in bioclimatic modeling today [6–8]. Moreover, the direct outputs of these systems are being transformed into derived climate data products on an unprecedented scale [9,10]. As a result, model tuning and variable selection, which are crucial aspects of any species distribution modeling effort, are becoming more complicated issues [11].

Part of the problem lies in the fact that MaxEnt, like many machine learning systems, acts on its inputs as a piece: predictors and observations must be memory-resident for the program to work [12]. This results in run-times and space requirements that scale linearly with the size of a model’s inputs. In most cases, these scaling properties pose few difficulties. But when the number of predictors under consideration becomes large, compute times can become impractically long, models can become overly complex, and efforts to understand any particular variable’s contribution to model formation, either as an aspect of model

analysis or as a way of selecting subsets of variables for further model refinement, can become challenging [11,13–16]. Clearly, an effective way of dealing with large environmental data sets that preserves the many advantages of MaxEnt while overcoming its current limitations would benefit the MaxEnt community.

In this study, we investigated the potential of a Monte Carlo method to help accomplish such an outcome. Monte Carlo optimizations are a common way of finding approximate answers to problems that are solvable in principle but lack a practical means of solution [17]. Our objective was to find a useful subset of predictors in a larger collection of environmental variables in a reasonable amount of time. Our proposed solution takes an ensemble approach wherein many MaxEnt runs, each drawing on a small random subset of variables, converges on a global estimate of the top contributing subset of variables in the larger collection.

Preliminary results suggest that the method reliably selects a subset of the original predictors that is capable of producing a well-tuned, parsimonious model of high quality. Since each model run is independent and uses a set number of variables, the method is totally parallelizable, independent of the scaling properties of MaxEnt, and amenable to implementation as an external memory algorithm. If proved to be effective, such an approach could provide a practical way of constructing MaxEnt models when there is a need to select a small set of predictors in a pool comprising a potentially very large number of predictors. This could lead to greater use of climate model outputs by the ecological research community and aid the search for viable predictors when variable selection through ecological reasoning is not apparent.

Materials and Methods

Cassin's Sparrow (*Peucaea cassinii* Woodhouse, 1852) is an elusive resident of arid shrub grasslands of Middle America and the Southwestern United States [18]. Desert-adapted birds, such as Cassin's Sparrow, appear to be especially vulnerable to climate change [19,20]. We chose Cassin's Sparrow as a target for our study as an example of a species whose study could benefit from the technical advances described here. Occurrence data was obtained from the Global Biodiversity Information Facility (GBIF) for the year 2016 [21]. After removing replicates, a total of 1865 records were acquired. To limit spatial extent and avoid pseudo-replication, we thinned the points to a radius of 16 km, which resulted in a total of 609 observations. For predictors, we used Worldclim's standard 19 Bioclimatic (bioclim) environmental variables at a resolution of 2.5 arc-minutes throughout (Table 1) [22]. We did not attempt to minimize collinearity by removing variables, because the current study focuses on an assessment of stochastic down-selection from a full variable set, and because MaxEnt has a demonstrated ability to account well for redundant variables [23].

We used MaxEnt Version 3.4.1 [24], R Version 4.0.1 [25], the ENMEval Version 0.3.0 R package [26], RStudio Version 1.2.5033 [27], and ENMTools Version 1.4.4 [28] running on a 2.8 GHz Intel Core i7 MacBook Pro with 16 GB of memory in the study. First, we developed a baseline model using the stand-alone MaxEnt program operated through its graphical user interface (GUI). MaxEnt users can apply various combinations of five mathematical transformations ('feature classes' or FCs) to predictor variables to enable more complex fits to the observational data. The available feature types for continuous variables are linear (L), quadratic (Q), hinge (H), product (P), and threshold (T) [1]. Users can also adjust a regularization multiplier (RM) to maximize predictive accuracy and offset the

Table 1. Worldclim Bioclimatic Variables.

bio01	Annual Mean Temperature
bio02	Mean Diurnal Range (Mean of monthly (max temp - min temp))
bio03	Isothermality (BIO2/BIO7) ($\times 100$)
bio04	Temperature Seasonality (standard deviation $\times 100$)
bio05	Max Temperature of Warmest Month
bio06	Min Temperature of Coldest Month
bio07	Temperature Annual Range (BIO5-BIO6)
bio08	Mean Temperature of Wettest Quarter
bio09	Mean Temperature of Driest Quarter
bio10	Mean Temperature of Warmest Quarter
bio11	Mean Temperature of Coldest Quarter
bio12	Annual Precipitation
bio13	Precipitation of Wettest Month
bio14	Precipitation of Driest Month
bio15	Precipitation Seasonality (Coefficient of Variation)
bio16	Precipitation of Wettest Quarter
bio17	Precipitation of Driest Quarter
bio18	Precipitation of Warmest Quarter
bio19	Precipitation of Coldest Quarter

overfitting that FC adjustments can introduce. We applied MaxEnt's default FC and RM settings (i.e. the "Auto features" setting) with 10 replicate cross-validation and jackknife evaluation of variable importance. By default, MaxEnt uses all feature classes and a regularization multiplier of 1.0 when there are more than 80 training samples, which was the case here [24]. Ten thousand background points were selected from across the study area following the recommendations of Phillips et al. [29] and Fourcade et al. [30]. We determined the average permutation importance for each variable in three replicated runs. The top six predictors in the three-run ensemble were used to develop the final MaxEnt baseline model.

We then developed an alternative method to select the top six variables that is based on random sampling. We implemented our Monte Carlo approach as an R script that invokes MaxEnt through ENMEval, which provides convenient control over model settings, built-in evaluation metrics, and improved performance [15,26]. To reduce variability and isolate outcomes as much as possible to the effects of the sampling process, we adopted a feature

class setting of LQHP and a regularization multiplier setting of 1.0 as fixed parameters in all the Monte Carlo runs. We defined ensemble, in this case, to mean a collection of 100 sprints, where each sprint consisted of ten model runs. A tally table was used to maintain a count of the number of times a variable was used in a model run along with a cumulative sum of the variable's permutation importance. The tally table thus provided the information needed to determine the average permutation importance of a predictor at any point along the way.

To process a sprint, we initialized each of its ten model runs with a random subset of environmental variables read from the filesystem. Random integers drawn from a uniform distribution ranging 1–19 corresponding to the 19 bioclim predictors were used to make the selection. At the conclusion of each model run, the tally table was updated appropriately. At the conclusion of each sprint, we computed a MaxEnt model using the six predictors in the original starting set having the highest average permutation importance values at that point. This process was repeated 100 times to produce a complete ensemble. We assessed the algorithm's performance in two ensembles. In the first, we chose two random variables for each sprint run; in the second, six random variables were used for each run. This resulted in an overall total of 2000 model runs.

The predictive distribution maps produced by the models were judged for reasonableness based on first-hand knowledge of the species, its habitat preferences, and known range [31]. We further compared model predictions to observational records from Cornell Lab's eBird citizen-scientist database [32]. We used the area under the operating curve (AUC) [33] as an indication of a model's classification accuracy (higher values indicating greater accuracy) and the Akaike information criterion corrected for small sample size (AICc) [34] as a measure of relative explanatory power (lower values indicating less

information loss). Model similarity was compared with Warren's I-statistic [35] and Schoener's D statistic [36] (higher values in both indicating greater similarity) using ENMTools. Single-processor run times were recorded to aid our understanding of algorithm performance and help identify opportunities for multiprocessor parallelization.

Results

On the basis of permutation importance, 13 of the 19 original bioclim variables were among the top ten most contributory predictors across all three replicated runs of the MaxEnt baseline: bio02, bio03, bio05, bio06, bio08–bio12, bio14, bio15, bio17, and bio18 (Table 2). Of those, bio02, bio05, and bio14 appeared in only one run each at 10th place. Bio18 showed strong dominance throughout. When performance was averaged across all three runs, the top six contributory variables in the ensemble collectively accounted for 65% of overall permutation importance (ensemble average). In descending order of importance, the top six predictors included bio18, bio03, bio10, bio15, bio11, and bio06. When these six top-contributing variables were used in a final MaxEnt run, the model's four most contributory variables (bio18, bio03, bio10, and bio15) accounted for approximately 86% of overall permutation importance, and its predicted habitat suitability distribution corresponded well with what is known about the natural history of the species and observational records for Cassin's Sparrow for the year 2016 (Fig 1) [32].

Table 2. Results of Maxent Baseline and Monte Carlo Selection Trials.

MODELS	Biodim Environmental Variables Permutation Importance														Permutation Importance % Top 3	Permutation Importance % Top 4	Total Run Count	Total Run Time		Avg Random Samples/Variable					
	bio01	bio02	bio03	bio04	bio05	bio06	bio07	bio08	bio09	bio10	bio11	bio12	bio13	bio14				bio15	bio16		bio17	bio18	bio19	AUC	AUCc
Maxent Baseline																									
Run #1	1.0	4.3	5.6	2.8	0.6	0.9	0.9	5.2	4.2	8.4	4.9	4.6	1.0	2.7	8.3	0.0	2.2	30.6	1.7			1	151	2.5	-
Run #2	1.7	2.9	6.4	1.9	3.5	9.9	0.5	3.5	3.9	5.4	7.7	4.8	1.0	3.5	8.3	0.8	0.5	30.1	3.6			1	120	2.0	-
Run #3	0.8	3.4	6.1	1.3	1.9	10.4	1.0	2.5	4.5	7.6	8.3	4.6	2.3	3.4	7.1	0.1	3.4	28.9	2.5			1	100	1.7	-
Ensemble avg	1.2	3.5	6.0	2.0	2.0	7.1	0.8	3.7	4.2	7.1	7.0	4.7	1.4	3.2	7.9	0.3	2.0	29.9	2.6						
Final model						6.6				13.8	7.3				10.9			47.4				1	18.0	0.3	-
Monte Carlo Selection																									
Ensemble #1	<i>(two random variables per sprint run)</i>																								
Sprint 025								10.8	9.3									14.3	25.9			250	181	3.0	26
Sprint 050								11.7	11.2									15.2	25.5			500	355	5.9	53
Sprint 100								6.2										4.1	39.7			1000	710	11.8	105
Ensemble #2	<i>(six random variables per sprint run)</i>																								
Sprint 025						3.6												3.8	43.9			250	438	7.3	79
Sprint 050						4.9												3.7	47.7			500	856	14.3	158
Sprint 100						5.5												2.6	42.4			1000	1793	29.9	316

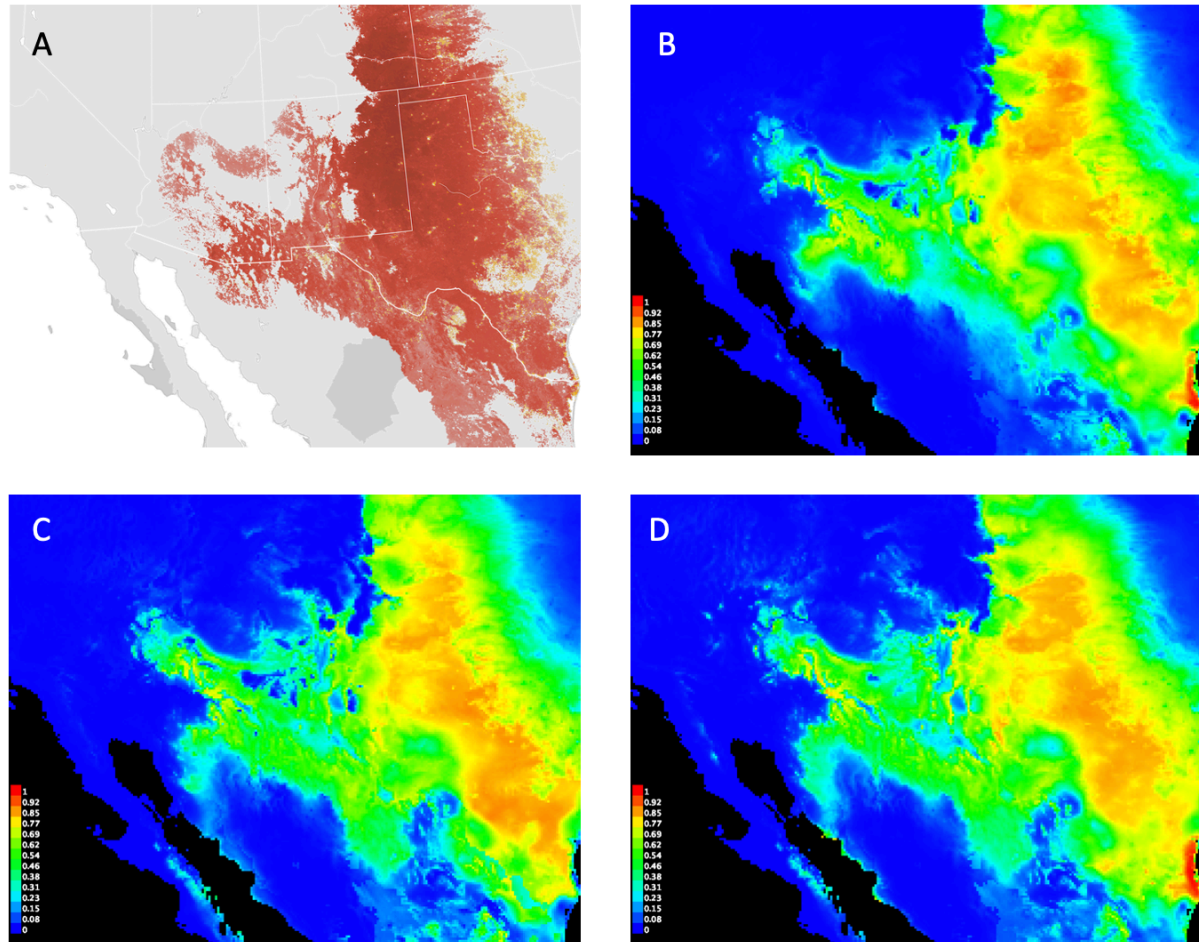


Fig 1. Cassin's Sparrow distribution maps. Cassin's Sparrow range map from Cornell Lab's eBird observational database (A). Predicted habitat suitability distributions from the MaxEnt baseline (B) and Monte Carlo Ensembles #1 (C) and #2 (D).

A distinct pattern of progression toward a stable subset of key variables was observed in the Monte Carlo ensembles (Figs 2,3). In both cases, the top three contributory variables among the top six were selected early in the sprint runs, and AICc values fluctuated within a narrow range around an average that changed little over the course of the selection process. Greater variability in the composition of the top six subset was seen in Ensemble #1 where two random variables at a time were selected for each sprint run (Table 2, Fig 2). In Ensemble #2, where six random variables at a time were selected for the MaxEnt runs, the top six variables were identified by the 25th sprint and had settled into their final rank order by sprint 54 (Fig 3). Ensemble #2 appeared to produce the best overall results and shared four variables in common with the top six selected by the MaxEnt baseline (bio03, bio06, bio11, and bio18) (Table 2). Ensemble #2's final model had the lowest overall AICc, and its four most contributory variables accounted for approximately 93% of overall permutation importance, the highest attained overall.

Ensemble #1 had only one variable in common with the top six selected by both the baseline run and Ensemble #2. What accounts for this difference is not immediately apparent; however, we speculate that the random pair-wise comparisons occurring in Ensemble #1 may alter the relative global influence of the collinearities known to exist in the bioclim variables [37–39]. The average number of times a variable was sampled appeared to have a marginal, positive influence on resulting model quality once an adequate minimum was attained. Ensemble #2 results suggest that at least 80 uniformly distributed samples per starting-set variable are needed to identify a reasonable top six set of variables; the best overall model resulted from over 300 samples per variable (Table 2).

Sprint	Runs	Time (min)	AICc	1st	2nd	3rd	4th	5th	6th	Sprint	Runs	Time (min)	AICc	1st	2nd	3rd	4th	5th	6th
1	10	7.58	12379.01	bio05	bio08	bio14	bio04	bio17	bio03	51	510	7.78	12233.56	bio18	bio16	bio13	bio08	bio05	bio09
2	20	8.04	12200.28	bio15	bio08	bio16	bio12	bio01	bio14	52	520	8.06	12259.47	bio18	bio16	bio13	bio08	bio05	bio09
3	30	6.99	12167.31	bio08	bio16	bio18	bio14	bio05	bio09	53	530	6.91	12241.89	bio18	bio16	bio13	bio08	bio05	bio09
4	40	6.87	12232.75	bio18	bio16	bio08	bio05	bio09	bio12	54	540	8.19	12247.57	bio18	bio16	bio13	bio08	bio05	bio09
5	50	7.31	12170.46	bio18	bio16	bio08	bio05	bio13	bio09	55	550	7.69	12226.76	bio18	bio16	bio13	bio08	bio05	bio09
6	60	7.17	12149.38	bio18	bio16	bio08	bio05	bio13	bio09	56	560	6.49	12210.05	bio18	bio16	bio13	bio08	bio05	bio09
7	70	6.69	12226.38	bio18	bio16	bio13	bio08	bio05	bio09	57	570	7.90	12205.59	bio18	bio16	bio13	bio08	bio05	bio09
8	80	7.73	12222.55	bio18	bio16	bio13	bio05	bio08	bio09	58	580	5.90	12254.37	bio18	bio16	bio13	bio08	bio05	bio09
9	90	6.40	12157.66	bio18	bio16	bio13	bio08	bio05	bio09	59	590	6.75	12273.07	bio18	bio16	bio13	bio08	bio05	bio09
10	100	7.89	12195.88	bio18	bio16	bio13	bio08	bio05	bio15	60	600	6.82	12236.51	bio18	bio16	bio13	bio08	bio05	bio09
11	110	8.48	12221.87	bio18	bio16	bio13	bio08	bio05	bio09	61	610	7.22	12243.03	bio18	bio16	bio13	bio08	bio05	bio09
12	120	7.11	12265.25	bio18	bio16	bio13	bio05	bio08	bio09	62	620	7.73	12198.01	bio18	bio16	bio13	bio08	bio05	bio09
13	130	5.60	12251.47	bio18	bio16	bio13	bio05	bio08	bio11	63	630	8.41	12338.98	bio18	bio16	bio13	bio08	bio05	bio09
14	140	8.10	12218.51	bio18	bio16	bio13	bio08	bio05	bio09	64	640	6.79	12245.66	bio18	bio16	bio13	bio08	bio05	bio09
15	150	7.16	12198.43	bio18	bio16	bio13	bio08	bio05	bio09	65	650	6.98	12251.55	bio18	bio16	bio13	bio08	bio05	bio09
16	160	7.32	12222.59	bio18	bio16	bio13	bio08	bio05	bio09	66	660	7.02	12348.13	bio18	bio16	bio13	bio08	bio05	bio09
17	170	6.99	12216.19	bio18	bio16	bio13	bio08	bio05	bio03	67	670	7.28	12239.33	bio18	bio16	bio13	bio08	bio05	bio09
18	180	6.95	12199.30	bio18	bio16	bio13	bio08	bio05	bio12	68	680	7.92	12343.54	bio18	bio16	bio13	bio08	bio05	bio09
19	190	6.20	12211.87	bio18	bio16	bio13	bio08	bio05	bio03	69	690	7.40	12237.15	bio18	bio16	bio13	bio08	bio05	bio09
20	200	8.30	12257.88	bio18	bio16	bio13	bio08	bio05	bio03	70	700	6.83	12222.24	bio18	bio16	bio13	bio08	bio05	bio09
21	210	7.08	12220.48	bio18	bio16	bio13	bio08	bio05	bio09	71	710	7.89	12280.63	bio18	bio16	bio13	bio08	bio05	bio09
22	220	6.94	12195.24	bio18	bio16	bio13	bio08	bio05	bio09	72	720	6.65	12234.48	bio18	bio16	bio13	bio08	bio05	bio09
23	230	7.53	12209.08	bio18	bio16	bio13	bio08	bio05	bio09	73	730	6.97	12256.76	bio18	bio16	bio13	bio08	bio05	bio09
24	240	8.17	12229.47	bio18	bio16	bio13	bio08	bio05	bio09	74	740	7.58	12279.34	bio18	bio16	bio13	bio08	bio05	bio09
25	250	6.61	12229.17	bio18	bio16	bio13	bio08	bio05	bio09	75	750	7.08	12210.23	bio18	bio16	bio13	bio08	bio05	bio09
26	260	7.67	12227.33	bio18	bio16	bio13	bio08	bio05	bio09	76	760	6.97	12232.46	bio18	bio16	bio13	bio08	bio05	bio09
27	270	7.68	12240.54	bio18	bio16	bio13	bio08	bio05	bio09	77	770	7.14	12298.37	bio18	bio16	bio13	bio08	bio05	bio09
28	280	5.29	12253.73	bio18	bio16	bio13	bio08	bio05	bio09	78	780	6.80	12220.16	bio18	bio16	bio13	bio08	bio05	bio09
29	290	7.69	12245.23	bio18	bio16	bio13	bio08	bio05	bio09	79	790	7.62	12214.21	bio18	bio16	bio13	bio08	bio05	bio09
30	300	8.49	12277.86	bio18	bio16	bio13	bio08	bio05	bio09	80	800	7.41	12227.38	bio18	bio16	bio13	bio08	bio05	bio09
31	310	6.38	12245.44	bio18	bio16	bio13	bio08	bio05	bio09	81	810	7.18	12250.13	bio18	bio16	bio13	bio08	bio05	bio09
32	320	7.53	12189.91	bio18	bio16	bio13	bio08	bio05	bio09	82	820	7.65	12264.45	bio18	bio16	bio13	bio08	bio05	bio03
33	330	7.08	12233.77	bio18	bio16	bio13	bio08	bio05	bio09	83	830	6.20	12245.57	bio18	bio16	bio13	bio08	bio05	bio03
34	340	6.79	12249.50	bio18	bio16	bio13	bio08	bio05	bio09	84	840	6.33	12216.66	bio18	bio16	bio13	bio08	bio05	bio03
35	350	6.52	12213.96	bio18	bio16	bio13	bio08	bio05	bio09	85	850	8.21	12242.54	bio18	bio16	bio13	bio08	bio05	bio03
36	360	6.80	12243.56	bio18	bio16	bio13	bio08	bio05	bio09	86	860	6.49	12282.52	bio18	bio16	bio13	bio08	bio05	bio03
37	370	7.71	12245.80	bio18	bio16	bio13	bio08	bio05	bio09	87	870	6.87	12256.19	bio18	bio16	bio13	bio08	bio05	bio03
38	380	7.96	12225.07	bio18	bio16	bio13	bio08	bio05	bio09	88	880	7.06	12251.20	bio18	bio16	bio13	bio08	bio05	bio03
39	390	6.42	12249.01	bio18	bio16	bio13	bio08	bio05	bio09	89	890	6.78	12237.27	bio18	bio16	bio13	bio08	bio05	bio03
40	400	5.70	12253.59	bio18	bio16	bio13	bio08	bio05	bio09	90	900	8.21	12278.73	bio18	bio16	bio13	bio08	bio05	bio09
41	410	7.08	12206.76	bio18	bio16	bio13	bio08	bio05	bio09	91	910	6.06	12307.16	bio18	bio16	bio13	bio08	bio05	bio09
42	420	7.06	12259.68	bio18	bio16	bio13	bio08	bio05	bio09	92	920	5.29	12210.69	bio18	bio16	bio13	bio08	bio05	bio09
43	430	6.41	12264.71	bio18	bio16	bio13	bio08	bio05	bio09	93	930	7.93	12245.61	bio18	bio16	bio13	bio08	bio05	bio03
44	440	6.70	12246.71	bio18	bio16	bio13	bio08	bio05	bio09	94	940	5.88	12229.42	bio18	bio16	bio13	bio08	bio05	bio03
45	450	6.06	12297.09	bio18	bio16	bio13	bio08	bio05	bio09	95	950	8.25	12236.33	bio18	bio16	bio13	bio08	bio05	bio03
46	460	6.93	12214.90	bio18	bio16	bio13	bio08	bio05	bio09	96	960	6.15	12280.85	bio18	bio16	bio13	bio08	bio05	bio03
47	470	6.81	12250.17	bio18	bio16	bio13	bio08	bio05	bio09	97	970	6.42	12273.77	bio18	bio16	bio13	bio08	bio05	bio03
48	480	6.71	12212.33	bio18	bio16	bio13	bio08	bio05	bio09	98	980	6.38	12247.16	bio18	bio16	bio13	bio08	bio05	bio03
49	490	6.95	12230.75	bio18	bio16	bio13	bio08	bio05	bio09	99	990	7.06	12214.92	bio18	bio16	bio13	bio08	bio05	bio03
50	500	7.06	12230.62	bio18	bio16	bio13	bio08	bio05	bio09	100	1000	7.17	12252.44	bio18	bio16	bio13	bio08	bio05	bio03

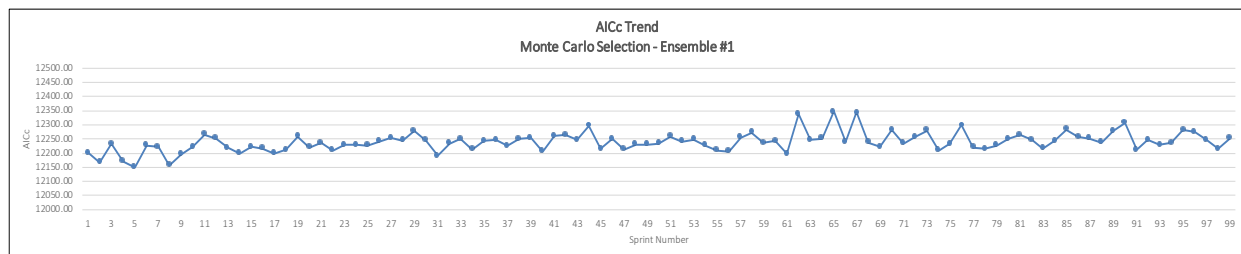


Fig 2. Monte Carlo Ensemble #1 results. Two random variables at a time were chosen for each MaxEnt sprint run. The sprint log on top shows the progressive selection of a stable set of top six variables in yellow. The graph on the bottom shows the narrow range of fluctuating AICc values over the course of the ensemble runs. Maximum and minimum AICc values are shown in red.

Sprint	Runs	Time (min)	AICc	1st	2nd	3rd	4th	5th	6th	Sprint	Runs	Time (min)	AICc	1st	2nd	3rd	4th	5th	6th
1	10	18.07	12185.35	bio18	bio06	bio12	bio16	bio03	bio11	51	510	18.53	12175.76	bio18	bio16	bio13	bio11	bio06	bio03
2	20	18.32	12182.82	bio18	bio16	bio06	bio12	bio03	bio14	52	520	20.15	12198.96	bio18	bio16	bio13	bio11	bio06	bio03
3	30	17.60	12232.96	bio18	bio16	bio13	bio06	bio03	bio15	53	530	20.48	12168.05	bio18	bio16	bio13	bio06	bio11	bio03
4	40	18.92	12180.38	bio18	bio16	bio13	bio03	bio12	bio14	54	540	17.22	12163.76	bio18	bio16	bio13	bio11	bio06	bio03
5	50	16.99	12207.91	bio18	bio16	bio13	bio03	bio06	bio04	55	550	17.28	12165.65	bio18	bio16	bio13	bio06	bio11	bio03
6	60	16.67	12191.40	bio18	bio16	bio13	bio03	bio06	bio11	56	560	20.14	12189.90	bio18	bio16	bio13	bio06	bio11	bio03
7	70	15.97	12167.45	bio18	bio16	bio13	bio06	bio03	bio11	57	570	20.86	12161.10	bio18	bio16	bio13	bio06	bio11	bio03
8	80	17.16	12188.02	bio18	bio16	bio13	bio06	bio03	bio11	58	580	19.10	12182.45	bio18	bio16	bio13	bio06	bio11	bio03
9	90	18.62	12161.62	bio18	bio16	bio13	bio06	bio03	bio11	59	590	17.63	12179.96	bio18	bio16	bio13	bio06	bio11	bio03
10	100	17.70	12179.81	bio18	bio16	bio13	bio06	bio11	bio03	60	600	19.88	12182.09	bio18	bio16	bio13	bio06	bio11	bio03
11	110	18.57	12163.00	bio18	bio16	bio13	bio06	bio11	bio03	61	610	18.92	12166.12	bio18	bio16	bio13	bio06	bio11	bio03
12	120	17.67	12161.86	bio18	bio16	bio13	bio06	bio11	bio03	62	620	19.17	12163.61	bio18	bio16	bio13	bio06	bio11	bio03
13	130	18.03	12195.77	bio18	bio16	bio13	bio06	bio11	bio03	63	630	18.74	12152.40	bio18	bio16	bio13	bio06	bio11	bio03
14	140	17.45	12162.09	bio18	bio16	bio13	bio06	bio11	bio03	64	640	20.34	12201.92	bio18	bio16	bio13	bio06	bio11	bio03
15	150	16.41	12153.21	bio18	bio16	bio13	bio06	bio11	bio03	65	650	19.29	12202.01	bio18	bio16	bio13	bio06	bio11	bio03
16	160	17.04	12169.11	bio18	bio16	bio13	bio06	bio11	bio03	66	660	19.49	12185.66	bio18	bio16	bio13	bio06	bio11	bio03
17	170	16.71	12173.60	bio18	bio16	bio13	bio06	bio11	bio03	67	670	19.05	12177.32	bio18	bio16	bio13	bio06	bio11	bio03
18	180	16.99	12160.81	bio18	bio16	bio13	bio06	bio11	bio03	68	680	18.96	12192.10	bio18	bio16	bio13	bio06	bio11	bio03
19	190	19.12	12165.07	bio18	bio16	bio13	bio06	bio11	bio03	69	690	18.21	12164.35	bio18	bio16	bio13	bio06	bio11	bio03
20	200	17.50	12160.62	bio18	bio16	bio13	bio06	bio11	bio03	70	700	20.04	12183.12	bio18	bio16	bio13	bio06	bio11	bio03
21	210	18.51	12151.23	bio18	bio16	bio13	bio06	bio11	bio03	71	710	18.20	12200.74	bio18	bio16	bio13	bio06	bio11	bio03
22	220	17.35	12281.16	bio18	bio16	bio13	bio06	bio11	bio03	72	720	19.34	12173.30	bio18	bio16	bio13	bio06	bio11	bio03
23	230	16.64	12200.43	bio18	bio16	bio13	bio06	bio11	bio03	73	730	17.71	12160.29	bio18	bio16	bio13	bio06	bio11	bio03
24	240	16.82	12172.95	bio18	bio16	bio13	bio06	bio11	bio03	74	740	19.39	12187.62	bio18	bio16	bio13	bio06	bio11	bio03
25	250	17.11	12176.06	bio18	bio16	bio13	bio06	bio11	bio03	75	750	17.44	12164.62	bio18	bio16	bio13	bio06	bio11	bio03
26	260	17.25	12166.65	bio18	bio16	bio13	bio06	bio11	bio03	76	760	17.59	12197.39	bio18	bio16	bio13	bio06	bio11	bio03
27	270	16.91	12203.32	bio18	bio16	bio13	bio06	bio11	bio03	77	770	17.78	12191.93	bio18	bio16	bio13	bio06	bio11	bio03
28	280	15.79	12182.27	bio18	bio16	bio13	bio06	bio11	bio03	78	780	19.23	12180.93	bio18	bio16	bio13	bio06	bio11	bio03
29	290	17.31	12186.27	bio18	bio16	bio13	bio06	bio11	bio03	79	790	18.43	12175.62	bio18	bio16	bio13	bio06	bio11	bio03
30	300	18.19	12238.20	bio18	bio16	bio13	bio06	bio11	bio03	80	800	18.13	12154.68	bio18	bio16	bio13	bio06	bio11	bio03
31	310	16.43	12177.40	bio18	bio16	bio13	bio06	bio11	bio03	81	810	18.11	12164.86	bio18	bio16	bio13	bio06	bio11	bio03
32	320	17.14	12159.92	bio18	bio16	bio13	bio06	bio11	bio03	82	820	20.02	12161.18	bio18	bio16	bio13	bio06	bio11	bio03
33	330	18.26	12149.59	bio18	bio16	bio13	bio06	bio11	bio03	83	830	19.08	12176.15	bio18	bio16	bio13	bio06	bio11	bio03
34	340	16.17	12178.16	bio18	bio16	bio13	bio06	bio11	bio03	84	840	19.17	12178.22	bio18	bio16	bio13	bio06	bio11	bio03
35	350	16.78	12153.80	bio18	bio16	bio13	bio11	bio06	bio03	85	850	20.01	12187.92	bio18	bio16	bio13	bio06	bio11	bio03
36	360	18.70	12158.82	bio18	bio16	bio13	bio06	bio11	bio03	86	860	17.05	12165.70	bio18	bio16	bio13	bio06	bio11	bio03
37	370	18.12	12179.92	bio18	bio16	bio13	bio11	bio06	bio03	87	870	17.78	12164.13	bio18	bio16	bio13	bio06	bio11	bio03
38	380	17.30	12187.95	bio18	bio16	bio13	bio11	bio06	bio03	88	880	18.00	12155.70	bio18	bio16	bio13	bio06	bio11	bio03
39	390	17.48	12208.89	bio18	bio16	bio13	bio11	bio06	bio03	89	890	16.87	12250.99	bio18	bio16	bio13	bio06	bio11	bio03
40	400	17.13	12203.25	bio18	bio16	bio13	bio11	bio06	bio03	90	900	17.98	12186.50	bio18	bio16	bio13	bio06	bio11	bio03
41	410	18.86	12167.30	bio18	bio16	bio13	bio11	bio06	bio03	91	910	17.64	12179.19	bio18	bio16	bio13	bio06	bio11	bio03
42	420	18.89	12166.73	bio18	bio16	bio13	bio11	bio06	bio03	92	920	17.89	12175.84	bio18	bio16	bio13	bio06	bio11	bio03
43	430	17.32	12170.09	bio18	bio16	bio13	bio11	bio06	bio03	93	930	15.97	12186.11	bio18	bio16	bio13	bio06	bio11	bio03
44	440	17.30	12181.02	bio18	bio16	bio13	bio11	bio06	bio03	94	940	17.25	12170.26	bio18	bio16	bio13	bio06	bio11	bio03
45	450	16.42	12194.59	bio18	bio16	bio13	bio11	bio06	bio03	95	950	16.83	12199.93	bio18	bio16	bio13	bio06	bio11	bio03
46	460	18.32	12148.92	bio18	bio16	bio13	bio11	bio06	bio03	96	960	16.43	12176.31	bio18	bio16	bio13	bio06	bio11	bio03
47	470	17.03	12196.84	bio18	bio16	bio13	bio11	bio06	bio03	97	970	18.13	12221.18	bio18	bio16	bio13	bio06	bio11	bio03
48	480	18.00	12178.66	bio18	bio16	bio13	bio11	bio06	bio03	98	980	17.58	12169.45	bio18	bio16	bio13	bio06	bio11	bio03
49	490	17.28	12168.32	bio18	bio16	bio13	bio11	bio06	bio03	99	990	15.28	12197.56	bio18	bio16	bio13	bio06	bio11	bio03
50	500	16.30	12198.48	bio18	bio16	bio13	bio11	bio06	bio03	100	1000	16.84	12151.75	bio18	bio16	bio13	bio06	bio11	bio03

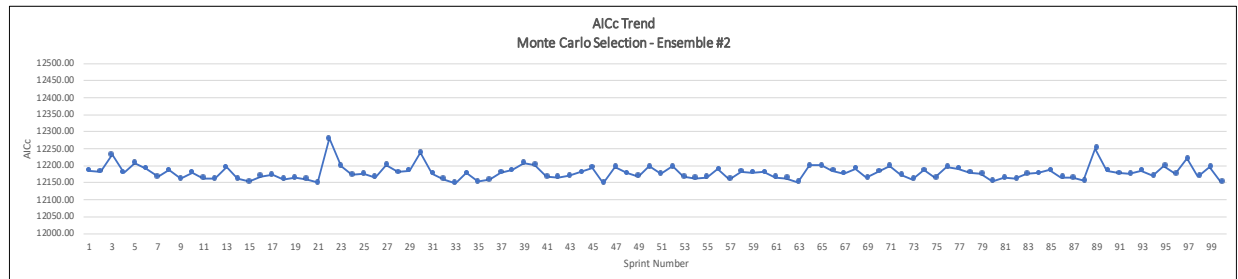


Fig 3. Monte Carlo Ensemble #2 results. Six random variables at a time were chosen for each MaxEnt sprint run. The sprint log on top shows the progressive selection of a stable set of top six variables in yellow. The graph on the bottom shows the narrow range of fluctuating AICc values over the course of the ensemble runs. Maximum and minimum AICc values are shown in red.

Discussion

The most striking outcome of the study is the similarity in results. Final models in the MaxEnt baseline and the Monte Carlo ensembles produced predicted habitat suitability distributions that are nearly indistinguishable from one another (Fig 1). Internal metrics likewise reveal little difference in outcomes, with AUC values ranging only from 0.801 to 0.818 and AICc ranging from 12,152 to 12,222, suggesting that the traditional MaxEnt runs and the Monte Carlo approach both produced reasonable models (Table 2). The two approaches also each identified four variables that collectively contributed more than 80% to the formulation of their respective models. Across the board, models showed a high degree of similarity in Schoener's D and the I-statistic (Table 3).

Table 3. Model Similarity Metrics.

<u>Schoener's D Statistic</u>										
MODELS ↓ →	Maxent-Run1	Maxent-Run2	Maxent-Run3	Maxent-Final	MC-E1-025	MC-E1-050	MC-E1-100	MC-E2-025	MC-E2-050	MC-E2-100
Maxent-Run1	1	0.9712	0.9738	0.9355	0.8629	0.8630	0.9044	0.8882	0.8903	0.8906
Maxent-Run2	x	1	0.9793	0.9397	0.8639	0.8648	0.9078	0.8915	0.8949	0.8947
Maxent-Run3	x	x	1	0.9354	0.8579	0.8586	0.9039	0.8871	0.8903	0.8902
Maxent-Final	x	x	x	1	0.8667	0.8673	0.9214	0.9104	0.9138	0.9142
MC-E1-025	x	x	x	x	1	0.9880	0.9006	0.8810	0.8801	0.8785
MC-E1-050	x	x	x	x	x	1	0.9026	0.8813	0.8804	0.8786
MC-E1-100	x	x	x	x	x	x	1	0.9393	0.9389	0.9365
MC-E2-025	x	x	x	x	x	x	x	1	0.9815	0.9810
MC-E2-050	x	x	x	x	x	x	x	x	1	0.9844
MC-E2-100	x	x	x	x	x	x	x	x	x	1

<u>Warren's I Statistic</u>										
MODELS ↓ →	Maxent-Run1	Maxent-Run2	Maxent-Run3	Maxent-Final	MC-E1-025	MC-E1-050	MC-E1-100	MC-E2-050	MC-E2-100	MC-E2-100
Maxent-Run1	1	0.9991	0.9994	0.9948	0.9760	0.9758	0.9874	0.9828	0.9833	0.9834
Maxent-Run2	x	1	0.9995	0.9949	0.9760	0.9760	0.9878	0.9831	0.9838	0.9839
Maxent-Run3	x	x	1	0.9945	0.9748	0.9748	0.9869	0.9823	0.9830	0.9831
Maxent-Final	x	x	x	1	0.9755	0.9753	0.9894	0.9871	0.9878	0.9880
MC-E1-025	x	x	x	x	1	0.9998	0.9887	0.9798	0.9797	0.9784
MC-E1-050	x	x	x	x	x	1	0.9889	0.9795	0.9794	0.9781
MC-E1-100	x	x	x	x	x	x	1	0.9910	0.9912	0.9906
MC-E2-025	x	x	x	x	x	x	x	1	0.9996	0.9994
MC-E2-050	x	x	x	x	x	x	x	x	1	0.9996
MC-E2-100	x	x	x	x	x	x	x	x	x	1

The most significant drawback identified in the study was the long run times.

MaxEnt's linear scaling behavior can be challenging in a single-processor environment. In the baseline runs, producing a single model through MaxEnt's GUI using our selected settings involved writing many files to disk and took from 18 minutes (with six variables) to over two hours (with all 19 variables). MaxEnt in the R environment outputs memory-resident objects, which results in faster run times. Still, with its repeated invocations of MaxEnt, Ensemble #2 took nearly 30 hours to complete (Table 2).

The Monte Carlo approach also exhibits linear scaling properties; however, its scaling behavior is not determined by the MaxEnt program, since each of the MaxEnt runs in the Monte Carlo method operates on a set number of predictors. The method's linear scaling property is determined, instead, by the need to adequately sample the starting set of environmental variables in order to obtain a good result. What makes a practical implementation of the Monte Carlo method possible is that each of its MaxEnt runs is entirely independent from all other runs in the ensemble. This high level of subtask independence is sometimes referred to as an “embarrassingly parallel” workload, which makes it relatively straightforward to implement in a cluster computing environment. If 1000 processors were recruited into service — which is becoming increasingly convenient with the proliferation of multiprocessor, high performance cloud computing — a 1000-run ensemble could conceivably take about as long as a single MaxEnt run.

The potential significance of this advantage becomes apparent when one considers the method's use with large collections of environmental data. The Monte Carlo approach described here provides an approximate solution to the problem of finding a useful k -size subset of an n -size collection of variables. In principal, there are $n! / [k!(n-k)!]$ variable

combinations to consider in such an evaluation, a staggering 27,000-plus six-variable subsets with the 19 bioclim variables alone. Algorithms that accomplish variable selection through stepwise removal or are otherwise bound to the linear scaling properties of underlying software components are inherently unable to exhaustively explore this combinatorial space. A Monte Carlo method makes such a search possible by randomly sampling the universe of possible combinations and returning approximate solutions in practical amounts of time, particularly if implemented as a high-performance cloud service.

The use of IPCC-class climate model outputs in efforts to assess the impacts of climate change on biodiversity and other ecosystem processes is growing. Exploring the potential of these massive data sets, expanded use of ensemble modeling, and the actual work of fitting models for the thousands of species scientists wish to study will require hundreds to thousands of projections [4,5]. An improved capacity to use large environmental data sets in MaxEnt modeling would benefit this work. We are encouraged to think that innovative use of Monte Carlo techniques might provide a helpful means of meeting this challenge.

Conclusions

This small-scale, proof-of-concept study leaves many practical and theoretical questions unanswered. Preliminary results, however, suggest that a Monte Carlo method could offer a viable means of selecting environmental predictors for MaxEnt models that is amenable to parallelization and scalable to large data sets, including externally-stored collections. This points to the possibility of near-real-time multiprocessor implementations that could enable broader and more exploratory use of global climate model outputs in environmental niche modeling and aid in the discovery of viable predictors. Next steps will focus on implementing this capability in NASA's Advanced Data Analytics Platform

(ADAPT) science cloud, evaluating the method's behavior using products generated by the Goddard Earth Observing System, Version 5 (GEOS-5) modeling system, extending stochasticity to feature class and regularization multiplier selection, developing automatic stopping rules, and evaluating the method's effectiveness in addressing research questions relating to climate change influences on Cassin's Sparrow abundance and distribution.

References

1. Phillips SJ, Anderson RP, Schapire RE. Maximum Entropy Modeling of Species Geographic Distributions. *Ecological Modelling*. 2006;190: 231–259. doi:10.1016/j.ecolmodel.2005.03.026.
2. Phillips SJ, Anderson RP, Dudík M, Schapire RE, Blair ME. Opening the black box: An open-source release of MaxEnt. *Ecography*. 2017;40: 887–893.
3. Elith J, Phillips SJ, Hastie T, Dudík M, Chee YE, Yates CJ. A statistical explanation of MaxEnt for ecologists. *Diversity and distributions*. 2011;17: 43–57.
4. Cavanagh RD, Murphy EJ, Bracegirdle TJ, Turner J, Knowland CA, Corney SP, et al. A Synergistic Approach for Evaluating Climate Model Output for Ecological Applications. *Frontiers in Marine Science*. 2017;4: 308. doi:10.3389/fmars.2017.00308
5. Araujo M, New M. Ensemble forecasting of species distributions. *Trends in Ecology & Evolution*. 2007;22: 42–47. doi:10.1016/j.tree.2006.09.010
6. Schnase JL. Climate Analytics as a Service. *Cloud Computing in Ocean and Atmospheric Sciences*. 2016. pp. 187–219. doi:10.1016/b978-0-12-803192-6.00011-6
7. Edwards PN. *A Vast Machine: Computer Models, Climate Data, and the Politics of Global Warming*. Cambridge, MA: MIT Press; 2010.
8. IPCC — Intergovernmental Panel on Climate Change. 2020 [cited 14 Mar 2020]. Available: <https://www.ipcc.ch/>
9. Harris RMB, Grose MR, Lee G, Bindoff NL, Porfirio LL, Fox-Hughes P. Climate projections for ecologists: Climate projections for ecologists. *Wiley Interdisciplinary Reviews: Climate Change*. 2014;5: 621–637. doi:10.1002/wcc.291
10. Responding to the Challenge of Climate and Environmental Change: NASA's Plan for a Climate-Centric Architecture for Earth Observations and Applications from Space. National Aeronautics and Space Administration; 2010. Available: <https://gmao.gsfc.nasa.gov/reanalysis/MERRA/>.

11. Araújo MB, Guisan A. Five (or so) challenges for species distribution modelling. *Journal of Biogeography*. 2006;33: 1677–1688. doi:10.1111/j.1365-2699.2006.01584.x
12. Duan Y, Edwards JS, Dwivedi YK. Artificial intelligence for decision making in the era of Big Data – evolution, challenges and research agenda. *International Journal of Information Management*. 2019;48: 63–71. doi:10.1016/j.ijinfomgt.2019.01.021
13. Galante PJ, Alade B, Muscarella R, Jansa SA, Goodman SM, Anderson RP. The challenge of modeling niches and distributions for data-poor species: a comprehensive approach to model complexity. *Ecography*. 2018;41: 726–736.
14. Zeng Y, Low BW, Yeo DCJ. Novel methods to select environmental variables in MaxEnt: A case study using invasive crayfish. *Ecological Modelling*. 2016;341: 5–13. doi:<https://doi.org/10.1016/j.ecolmodel.2016.09.019>
15. Muscarella R, Galante PJ, Soley-Guardia M, Boria RA, Kass JM, Uriarte M, et al. ENMeval: An R package for conducting spatially independent evaluations and estimating optimal model complexity for MaxEnt ecological niche models. *Methods in Ecology and Evolution*. 2014;5: 1198–1205. doi:10.1111/2041-210X.12261
16. Radosavljevic A, Anderson RP. Making better MaxEnt models of species distributions: complexity, overfitting and evaluation. Araújo M, editor. *Journal of Biogeography*. 2014;41: 629–643. doi:10.1111/jbi.12227
17. Kroese DP, Brereton T, Taimre T, Botev ZI. Why the Monte Carlo method is so important today: Why the MCM is so important today. *Wiley Interdisciplinary Reviews: Computational Statistics*. 2014;6: 386–392. doi:10.1002/wics.1314
18. Dunning, Jr. JB, Bowers, Jr. RK, Suter SJ, Bock CE. Cassin’s Sparrow (*Peucaea cassinii*), Version 1.0. In: *Birds of the World* (P. G. Rodewald, Editor) [Internet]. 2020 [cited 22 May 2020]. Available: <https://doi.org/10.2173/bow.casspa.01>
19. Iknayan KJ, Beissinger SR. Collapse of a desert bird community over the past century driven by climate change. *Proc Natl Acad Sci USA*. 2018;115: 8597. doi:10.1073/pnas.1805123115
20. Radchuk V, Reed T, Teplitsky C, van de Pol M, Charmantier A, Hassall C, et al. Adaptive responses of animals to climate change are most likely insufficient. *Nature Communications*. 2019;10: 3109. doi:10.1038/s41467-019-10924-4
21. GBIF. 2020 [cited 22 May 2020]. Available: <https://www.gbif.org/>
22. Bioclimatic variables — WorldClim 1 documentation. 2020 [cited 22 May 2020]. Available: <https://worldclim.org/data/bioclim.html>
23. Feng X, Park DS, Liang Y, Pandey R, Papeş M. Collinearity in ecological niche modeling: Confusions and challenges. *Ecology and Evolution*. 2019;9: 10365–10376. doi:10.1002/ece3.5555

24. MaxEnt Version 3.4.1 Download Site. [cited 22 May 2020]. Available: https://biodiversityinformatics.amnh.org/open_source/MaxEnt/
25. R: The R Project for Statistical Computing. [cited 22 May 2020]. Available: <https://www.r-project.org/>
26. Muscarella R, Galante PJ, Soley-Guardia M, Boria RA, Kass JM, Anderson MU and RP. ENMeval: Automated Runs and Evaluations of Ecological Niche Models. 2018. Available: <https://CRAN.R-project.org/package=ENMeval>
27. RStudio | Open source & professional software for data science teams. [cited 27 May 2020]. Available: <https://rstudio.com/>
28. Warren DL, Glor RE, Turelli M. ENMTools: a toolbox for comparative studies of environmental niche models. *Ecography*. 2010 [cited 27 Mar 2020]. doi:10.1111/j.1600-0587.2009.06142.x
29. Phillips SJ, Dudík M, Elith J, Graham CH, Lehmann A, Leathwick J, et al. Sample selection bias and presence-only distribution models: implications for background and pseudo-absence data. *Ecological applications*. 2009;19: 181–197.
30. Fourcade Y, Engler JO, Rödder D, Secondi J. Mapping Species Distributions with MAXENT Using a Geographically Biased Sample of Presence Data: A Performance Assessment of Methods for Correcting Sampling Bias. Valentine JF, editor. *PLoS ONE*. 2014;9: e97122. doi:10.1371/journal.pone.0097122
31. Schnase JL, Grant WE, Maxwell TC, Leggett JJ. Time and energy budgets of Cassin's sparrow (*Aimophila cassinii*) during the breeding season: evaluation through modelling. *Ecological Modelling*. 1991;55: 285–319.
32. (*Peucaea cassinii*) - Species Map - eBird. 2020 [cited 31 May 2020]. Available: <https://ebird.org/map/casspa>
33. Fielding AH, Bell JF. A review of methods for the assessment of prediction errors in conservation presence/absence models. *Environmental Conservation*. 1997;24: 38–49. doi:10.1017/S0376892997000088
34. H. Akaike. A new look at the statistical model identification. *IEEE Transactions on Automatic Control*. 1974;19: 716–723. doi:10.1109/TAC.1974.1100705
35. Warren DL, Glor RE, Turelli M. Environmental niche equivalency versus conservatism: quantitative approaches to niche evolution. *Evolution*. 2008;62: 2868–2883. doi:10.1111/j.1558-5646.2008.00482.x
36. Schoener TW. The *Anolis* Lizards of Bimini: Resource Partitioning in a Complex Fauna. *Ecology*. 1968;49: 704–726. doi:10.2307/1935534

37. Tang Y, Winkler JA, Viña A, Liu J, Zhang Y, Zhang X, et al. Pearson pairwise correlation matrix between the bioclimatic variables. 2018. doi:10.1371/journal.pone.0189496.g003
38. O'Donnell MS, Ignizio D a. Bioclimatic Predictors for Supporting Ecological Applications in the Conterminous United States. Reston, VA: US Geological Survey; 2012 p. 10. Report No.: 691. Available: <https://pubs.usgs.gov/ds/691/>
39. Dormann CF, Elith J, Bacher S, Buchmann C, Carl G, Carré G, et al. Collinearity: a review of methods to deal with it and a simulation study evaluating their performance. *Ecography*. 2013;36: 27–46. doi:10.1111/j.1600-0587.2012.07348.x