

1 **Robustness and accuracy of Maxent niche modelling for *Lactuca* species**  
2 **distributions in light of collecting expeditions**

3

4 **short title: Maxent model and Lactuca collecting expedition**

5

6 **MMP Cobben\*<sup>1,4</sup>, R van Treuren<sup>1</sup>, NP Castañeda-Álvarez<sup>2</sup>, CK Khoury<sup>2,3</sup>, C Kik<sup>1</sup>, TJJ**  
7 **van Hintum<sup>1</sup>**

8

9 \*corresponding author, e-mail: m.cobben@nioo.knaw.nl

10 <sup>1</sup> Centre for Genetic Resources, The Netherlands (CGN)

11 P.O. Box 16

12 6700 AA Wageningen

13 The Netherlands

14 <sup>2</sup> International Center for Tropical Agriculture (CIAT)

15 Apartado Aéreo 6713

16 Calí, Colombia

17 <sup>3</sup> Centre for Crop Systems Analysis, Wageningen University

18 P.O. Box 16

19 6700 AA Wageningen

20 The Netherlands

21 <sup>4</sup> Netherlands Institute of Ecology (NIOO-KNAW)

22 P.O. Box 50

23 6700 AB Wageningen

24 The Netherlands

25 **ABSTRACT**

26 Niche modelling software can be used to assess the probability of detecting a population of a  
27 plant species at a certain location. Here we used the distribution of wild relatives of lettuce  
28 (*Lactuca spp.*) to investigate the applicability of Maxent species distribution models for  
29 collecting missions. Geographic origin data of genebank and herbarium specimens and  
30 climatic data of the origin locations were used as input. For *L. saligna* we varied the input  
31 data by omitting specimens from different parts of the known distribution area to assess the  
32 robustness of the predicted distributions. Furthermore, we examined the accuracy of the  
33 modelling by comparing the predicted probabilities of population presence against recent  
34 expedition data for the endemic *L. georgica* and the cosmopolitan *L. serriola*. We found  
35 Maxent to be quite robust in its predictions, although its usefulness was higher for endemic  
36 taxa compared to more widespread species. The exclusion of occurrence data from the  
37 perceived range margins of the species can result in important information about local  
38 adaptation to distinct climatic conditions. We discuss the potential for enhanced use of  
39 Maxent in germplasm collecting planning.

40

41 Key words: Genebank, collecting expedition, *Lactuca*, species distribution modelling, niche  
42 modelling, Maxent, plant genetic resources

43

44

45 **INTRODUCTION**

46 Species distribution models have been used to predict distributions of a wide range of taxa  
47 (e.g. Guisan and Thuiller, 2005; Araújo and Guisan, 2006), and are increasingly employed for  
48 crop wild relatives (Guarino, 1995; Afonin and Greene, 1999; Greene *et al.*, 1999; Jarvis *et*  
49 *al.*, 2005; Maxted *et al.*, 2008; Parra-Quijano *et al.*, 2012). The Crop Wild Relatives &  
50 Climate Change project of the Global Crop Diversity Trust and Millennium Seed Bank  
51 Partnership, Royal Botanic Garden, Kew ([www.cwrdiversity.org](http://www.cwrdiversity.org)) uses a gap analysis  
52 methodology developed by Ramírez-Villegas *et al.* (2010), which incorporates the use of the  
53 maximum entropy model "Maxent" (Phillips *et al.*, 2006) to support collection planning for  
54 crop wild relatives. The input data for Maxent include the geographic origins of both  
55 genebank and herbarium specimens, and statistics representing the current climate, i.e. a set of  
56 temperature and precipitation parameters. The Maxent output distribution maps intend to give  
57 an indication of locations where the species may be present.

58 Collectors of plant genetic resources (PGR) are interested in material with new genetic  
59 diversity, preferably from species that can be crossed with the cultivated taxa for crop  
60 improvement. Such PGR can typically be collected from regions that have not been sampled  
61 previously. The gap-analysis method for crop genepools (Ramírez-Villegas *et al.*, 2010)  
62 assesses the priority with which a certain crop wild relative should be collected using the  
63 sampling representativeness score (SRS), estimating gross representation in genebanks; and  
64 the geographic representativeness score (GRS) and the environmental representativeness  
65 score (ERS), of which the GRS estimates comprehensiveness of genebank collections  
66 regarding the geographic range of the species, and ERS estimates comprehensiveness based  
67 on a principal component analysis (PCA) of the full environmental range of the modelled  
68 species. As such, it is possible to get an overview of populations that are underrepresented in  
69 genebanks and which may contain novel genetic diversity. The GRS and ERS depend upon

70 Maxent distribution models. In the method, it is implicitly assumed that the herbaria data  
71 provide full coverage of the distribution area of the species. In practice, herbaria data are  
72 incomplete and sampling bias can result in the systematic exclusion of certain regions of the  
73 species distribution from presence data. As a result, the distribution deduced with an  
74 associative species distribution model such as Maxent is not necessarily the complete  
75 distribution of the species. Maxent improves its models by excluding part of the presence data  
76 from the training sample to subsequently use it as test sample. However, this test sample is  
77 selected randomly, so does not systematically exclude a specific area, which mimicks the  
78 detection of an undiscovered region of the species distribution.

79         Here we use the distribution of wild relatives of lettuce (*Lactuca spp.*) as a case study  
80 to investigate how the predicted distribution of *L. saligna* depends on the input occurrence  
81 data. In addition, we compare the Maxent distribution predictions against expedition data of  
82 the endemic *L. georgica* and the cosmopolitan *L. serriola*. The results are utilized to discuss  
83 the applicability of Maxent to support PGR collecting missions.

84

85

## 86 MATERIAL AND METHODS

### 87 Study material

88 The wild relatives of lettuce were chosen as the focus of the case study due to high available  
89 knowledge and data for this crop, and the coincidence with a *Lactuca* collecting mission in  
90 the Trans Caucasus organised by the Centre for Genetic Resources, the Netherlands (CGN) in  
91 2013. Moreover, lettuce relatives represent a wide variety of both niche and distribution sizes,  
92 ranging from pan-temperate distributions for *L. serriola* (D'Andrea *et al.*, 2009; Alexander,  
93 2013) and *L. saligna*, to the narrow endemism of *L. georgica* in the Caucasus region (Zohary,  
94 1991) and thus permit an evaluation of modelling methods for a variety of species types.

95 We collected information about known *Lactuca* populations from both herbarium and  
96 genebank databases (Table 1). The International Lactuca Database (ILDB), Eurisco and the  
97 National Plant Germplasm System (NPGS) were used as the sources of information about  
98 genebank accessions. For herbarium samples we consulted the Global Biodiversity  
99 Information Facility (GBIF) as well as obtained occurrence records directly from herbaria and  
100 researchers (Supplement Table S1). We used the inventory by Van Treuren *et al.* (2012) to  
101 select the species known to belong to the genepool of cultivated lettuce (*L. sativa*) (Table 1)  
102 and to check for synonyms (Supplement Table S2) and reassignments from other genera. Van  
103 Treuren *et al.* (2012) did a survey of the International Plant Names Index, which revealed a  
104 total number of 538 *Lactuca* species, of which 357 referred to synonyms and basionyms,  
105 whereas for another 51 the taxonomic status and their belonging to the genus *Lactuca* was  
106 questionable. Of the remaining 130 species, 20 are generally considered to be part of the  
107 lettuce gene pool (Table 1).

108 The data were cleaned, records without and with only coarse geographic information  
109 were removed. In addition, we removed the duplicate species-specific locations. For the

110 remaining locations we used DIVA-GIS (vs 7.5, [www.diva-gis.org](http://www.diva-gis.org)) (Hijmans et al 2012) to  
111 cross-check the match between longitude/latitude combination and the stated country.

112

### 113 **Species distribution modelling**

114 Current climatic conditions were downloaded from [www.worldclim.org](http://www.worldclim.org) at a scale of 2.5 arc  
115 minutes, including 19 bioclimatic variables (Table 2). Maxent (Phillips *et al.*, 2006) uses a  
116 presence-only dataset as inputs and background points as pseudo-absences. The required  
117 10,000 background points were selected randomly in each of the continents with species-  
118 specific occurrence data to ameliorate sample bias (VanDerWal *et al.*, 2009). Populations  
119 which were located outside the grid of climate cells, e.g. in the sea, were removed from the  
120 dataset.

121

### 122 **Robustness of predicted distributions**

123 *L. saligna* was chosen to assess the robustness of the Maxent projections in the marginal areas  
124 of the distribution range, since this species has a pan-temperate distribution with sufficient  
125 data points to exclude specific regions. Together with *L. serriola* and *L. virosa*, *L. saligna*  
126 serves as an important source of novel diversity for exploitation in the development of novel  
127 lettuce varieties by breeding companies (Lebeda *et al.*, 2009).

128 Occurrence locations of *L. saligna* included North America, Australia, and Europe and from  
129 there extending into the eastern Mediterranean and Caucasus (Supplement Figure S1). To  
130 assess the robustness of the projections, we excluded the *L. saligna* occurrence samples in,  
131 respectively, the Greek region (GRC, 88 occurrences), the Israeli region (ISR, 32  
132 occurrences) and the Eurasian region (EUR, 684 occurrences) (Supplement Figure S1). The  
133 Eurasian region was excluded to serve as a benchmark for the other predictions. The Greek  
134 and Israeli regions are very distinct marginal areas within the Eurasian distribution, from a

135 geographic point of view (Supplement Figure S1). Maxent was run for each of the three new  
136 datasets and the output was compared with the model based on the original, complete dataset.  
137 We used a ten fold division of the input data, each fold replicating the model using the  
138 consecutive parts as the test sample, while the remaining 90% of the input data was used as  
139 the training sample. For visual comparison we used the Maxent projections, based on this  
140 tenfold cross-validation for each dataset. In addition, we correlated the Maxent estimated  
141 probabilities of occurrence of each 2.5 arc minute cell in each of the excluded regions that  
142 resulted from the models with the reduced and the original dataset. This indicates how the  
143 estimated probability of occurrence in a single cell changes when the occurrence data is  
144 excluded from the dataset on which the probability model was based.

145

#### 146 **Relationship with expedition data**

147 A *Lactuca* expedition to the Trans Caucasus was organised by CGN in 2013, which resulted  
148 in 94 unique collection locations. We compared the presence and absence locations of *L.*  
149 *georgica* and *L. serriola* with the model predictions for these locations. The model predictions  
150 were made *without including* the data of the sites visited during the expedition. So, while the  
151 modelling was done after the completion of the expedition, we mimicked the data availability  
152 prior to the collecting expedition, as if the modelling was done in its preparation. *L. georgica*  
153 and *L. serriola* were chosen because these 2 species were collected in a fair number of  
154 populations (32 and 55, respectively). In addition, they represent two opposites of the  
155 endemism spectrum, with *L. georgica* being endemic to the Trans Caucasus, while *L. serriola*  
156 has a pan-temperate distribution, with very many data points (Table 1). To investigate the  
157 effect of zooming in on the target area, a new *L. serriola* model was made using only the 133  
158 known occurrence samples from the expedition region, again without including data from the  
159 newly sampled populations. The predicted probability of occurrence was determined for the

160 2.5 arc minute grid cell in which a population was found. All probabilities were classified in  
161 categories ranging from 0.1 to 1.0, and the number of presence and absence locations were  
162 summed per class.

163



## 164 **RESULTS**

165

### 166 **Quality of the models**

167 The four Maxent models for *L. saligna* (the full dataset and those excluding respectively  
168 GRC, ISR and EUR) and the model for *L. georgica* classified as valid models (Supplement  
169 Table S3) according to the gap-analysis protocol (Ramírez-Villegas *et al.*, 2010), for which  
170 the average test AUC (area under the curve) should be larger than 0.7, the standard deviation  
171 of the test AUC smaller than 0.15 and the proportion of the potential distribution area with a  
172 standard deviation of the estimated probability of occurrence  $> 0.15$  should be smaller than  
173 10%. The *L. serriola* model including all occurrence data was not considered valid, as a result  
174 of the low average test AUC of 0.65. However, when the analysis was limited to the  
175 expedition region this resulted in a valid model for *L. serriola*.

176

### 177 **Robustness of model predictions for *L. saligna***

178 For *L. saligna*, excluding the occurrence data points in the Greek region of investigation  
179 (GRC) resulted in small changes in the estimated probabilities of occurrence in this region.  
180 The patterns in probability distributions were very similar between the two models (Fig. 1  
181 maps) with a good correlation between both models (Fig. 1 scatter plot) with a small decrease  
182 in estimated probabilities in the model where the occurrences in GRC were omitted as  
183 compared to the model that included all occurrences.

184 Excluding the occurrence data points in the Israeli region of investigation (ISR) led to  
185 large changes in the estimated probabilities of occurrence of *L. saligna* in this region (Fig. 2).  
186 The similarity in the patterns of probability distributions was recognisable (Fig. 2 maps), but  
187 the decrease in estimated probabilities from the model in which the local occurrence data  
188 were omitted was substantial (Fig. 2 scatter plot).

189           When all the occurrences in the Eurasian region of investigation (EUR) were  
190 excluded, the probabilities of occurrence of *L. saligna* changed drastically (Fig. 3 scatter plot),  
191 showing increases and decreases depending on the location. The pattern of the potential  
192 distribution changed from an emphasis on Western Europe to a most probable occurrence in  
193 the Middle-East and Central Asia (Fig. 3 maps). The maximum estimated probability of  
194 occurrence increased from 0.75 to 0.97. However, the margins of the potential distribution of  
195 *L. saligna* were very similar between both models.

196

### 197 **Model predictions in relation to expedition data**

198 Fig. 4 shows the presence and absence of both *L. georgica* and *L. serriola* in 94 unique  
199 locations sampled during the CGN *Lactuca* expedition in 2013 in relation to the Maxent  
200 estimated probabilities of occurrence in the matching 2.5 arc minute grid cells. The model  
201 performed quite well for *L. georgica*, as at locations with low estimated probabilities of  
202 occurrence only absence of the species was observed, while at the locations with high  
203 estimated probabilities of occurrence observed presence was considerably higher than absence  
204 of the species . In contrast, the *L. serriola* projection showed very little differentiation across  
205 the expedition region, restricting the estimated probability of occurrence to only a few classes  
206 in the middle of the potential range. The majority of locations fell within the 0.4-0.5  
207 probability class, at which more or less equal numbers of absences and presences were  
208 observed. The new *L. serriola* model, excluding all occurrence data but the ones in the  
209 expedition region, showed much more differentiation. However, the estimated probability of  
210 occurrence appeared a poor predictor for the presence and absence of the species.

211

212

213 **DISCUSSION**

214

215 **Species distribution modelling**

216 Species distribution models have been used for a number of decades and for many purposes  
217 (review by Elith and Leathwick, 2009), such as the consequences of climate change on  
218 species distributions and the assessment of the distribution of an invasive species. With the  
219 growing availability of digital records from natural history museums, herbaria, and genebanks  
220 coupled with the demand for mapped predictions, the incentive to put this source of presence-  
221 only data to use has been increasing (Elith and Leathwick, 2009). Many different modelling  
222 techniques exist and the discussion on how to best model presence-only data continues. It is  
223 now common to compare presence data with background or pseudo-absence data, by e.g.  
224 using regression methods such as generalised linear models (GLM), generalised additive  
225 models (GAM) or multivariate adaptive regression splines (MARS), but also GARP (Genetic  
226 Algorithm for Rule Set Production, Stockwell and Peters, 1999), ENFA (Ecological Niche  
227 Factor Analysis, Hirzel *et al.*, 2002) and Maxent (Phillips *et al.*, 2006), in chronological order.  
228 Elith *et al.* (2006) review and compare these methods. More recently, a platform was  
229 developed to combine different techniques for ensemble forecasting of species distributions  
230 (Araújo and New, 2007), called BIOMOD (R-package, <http://R-Forge.R-project.org>) (Thuiller  
231 *et al.*, 2009). This software is able to fit and compare different model classes in an attempt to  
232 reach more robust forecasts by treating the methodological uncertainties in different models.  
233 The latest version of BIOMOD, biomod2, also includes Maxent as one of the techniques.

234         In this study we have chosen to use Maxent as the modelling technique to align our  
235 results with the gap analysis method (Ramírez-Villegas *et al.*, 2010), which is used as input  
236 for many planned collecting expeditions in the Crop Wild Relatives & Climate Change  
237 project. With this project in mind, we have aimed to provide information about Maxent's

238 robustness and accuracy with regard to such collecting expeditions. Apart from these  
239 considerations, the accessibility and relative ease of use of the Maxent software compared to  
240 the others mentioned are valuable assets for application by non-experts in the planning of  
241 collecting expeditions. However, for increasing the strength of the predictions, using and  
242 comparing different modelling techniques as is done in BIOMOD is likely a valuable  
243 contribution. Compared to Maxent, this does require increased need for processing power,  
244 technical knowledge and time, which might limit the usability of BIOMOD for the non-  
245 expert.

246

#### 247 **Robustness of *L. saligna* models**

248 From the comparison of the different *L. saligna* models, we conclude that the Greek region is  
249 not climatically distinct within the known *L. saligna* distribution area (Fig. 1). In contrast, the  
250 Israeli region appeared very distinct (Fig. 2). Comparing the response curves of the Maxent  
251 predictions to the different environmental variables (Fig. S2), we found the Israeli region to  
252 differ from other areas within the distribution range especially regarding the mean annual  
253 temperature (BIO1), the maximum temperature of the warmest month (BIO5), the variation in  
254 the precipitation over the seasons (BIO15), and the total precipitation of the warmest quarter  
255 (BIO18). Fig. 2 indicates that the model is not capable of predicting this latter region as a  
256 potential distribution area when data from this region are excluded, while the unique climatic  
257 conditions may indicate the presence of potentially interesting diversity. For the Eurasian  
258 region of investigation, although the estimated local probabilities differ substantially between  
259 the two models (Fig. 3), the distribution borders are very similar. The changes in distribution  
260 pattern indicate that the south-eastern region, where the probabilities in occurrence increase  
261 when the Eurasian occurrences are excluded, is climatically more similar to the other regions  
262 in the world where *L. saligna* is sampled. From this it follows that, in addition to the Israeli

263 region, Western Europe may be considered a relatively exceptional climatic region in the *L.*  
264 *saligna* distribution area.

265 From this analysis we conclude that Maxent's outputs are generally robust, yet on a local  
266 scale, on which a collecting mission is typically planned, estimated probabilities of  
267 occurrence can differ to a larger extent. This depends on whether populations in the specific  
268 region have been sampled before (and the data used as input for the model), and on the  
269 climatic relatedness to other sampled regions within the distribution. It is important to note  
270 that all species distribution models project suitable habitat by climate association, and thus  
271 that none can predict potential distribution in regions that are climatically unique compared to  
272 regions where the species has been sampled. Particularly for the purpose investigated here, in  
273 search of unique plant genetic resources, this is an important limitation. However, excluding  
274 occurrence data from the perceived range margins of the species may result in important  
275 information about local adaptation to distinct climatic conditions. A principal component  
276 analysis of the climatic data from the occurrence locations can provide the same information  
277 and may be a good starting point for such an analysis. There are other possibilities to get  
278 information about local adaptation, e.g using SNP (single nucleotide polymorphism) data to  
279 confirm that a population is genetically different from similar populations, or checking  
280 phenotypic characteristics. Such methods tend to be used at a later stage, while the method we  
281 suggest here can be done with currently available information about population locations to  
282 get an indication about where such useful locally adapted populations may be present.

283 *L. saligna* is native to Eurasia and North-Africa (GRIN, 2014; Lebeda *et al.*, 2004a;  
284 Lebeda *et al.*, 2004b), and when omitting the Eurasian occurrence locations, we have  
285 essentially predicted the native distribution of the species from its non-native distribution.  
286 While this is not a logical procedure for collection planning, it gives indication of robustness  
287 of the predictions and the possibility to predict species' presence in a region where the species

288 has not been sampled. Omitting many data points from the native area of the species then  
289 makes an interesting benchmark study with which to compare the other results from the Greek  
290 and Israeli regions. Interestingly, this procedure could be used to indicate the potential origin  
291 of populations in the non-native distribution area (Alexander, 2013).

292

### 293 **Relationship with expedition data**

294 *L. georgica* is an endemic (Zohary, 1991) species, that lives in an equilibrium environment  
295 where natural competition determines the distribution of species. *L. serriola* is a ruderal  
296 species, thus living in disturbed environments (Grime, 1977), with a cosmopolitan distribution  
297 (D'Andrea *et al.*, 2009; Alexander, 2013). The latter's widespread distribution and the very  
298 many data points that are included in the initial model (Table 1) likely account for the rather  
299 undifferentiated model projections observed for *L. serriola*. It is interesting to note that,  
300 although not very informative, this model does provide a good prediction of its occurrence,  
301 with a similar number of presence and absence locations across the 50% probability region.  
302 When we limited the model to the expedition region, the number of probability classes  
303 increased substantially. However, actual presence and absence locations appeared to correlate  
304 poorly with the corresponding probability of occurrence according to the modelling (Fig. 4c).  
305 The *L. serriola* absence locations are mostly *L. georgica* presence locations, representing  
306 fairly undisturbed habitats. In addition, the ruderal nature of *L. serriola*, combined with its  
307 global distribution, explains its relative insensitivity to climatic conditions. Thus, the Maxent  
308 model for this global, ruderal species would not have been informative or otherwise useful for  
309 the Trans Caucasus expedition, not even when the projections would have been restricted to  
310 the region of interest. This is in line with the gap analysis protocol as suggested by Ramírez-  
311 Villegas *et al.* (2010), who excluded weedy species from the analysis. Here it needs to be  
312 noted that the collection of wide-spread and ruderal species such as *L. serriola* does not

313 require any modelling support, since such species can be easily located. In contrast, expected  
314 and actual presence data correlated well for the endemic *L. georgica*, living in pristine habitat.  
315 In fact, a population was sampled in a region where the local experts had not expected it,  
316 while the model predicted a high probability of occurrence at this location. In the case that the  
317 projection for *L. georgica* would have been available prior to the collecting mission, the  
318 expedition route would have been slightly adjusted to explore a nearby region where also high  
319 probabilities of occurrence were estimated. Thus, Maxent distribution models may be useful  
320 to support collecting missions and based on our findings this may particularly apply to  
321 endemic species growing in relatively undisturbed habitats.

322

### 323 **Concluding remarks**

324 Based on our results, we recommend organisers of collecting missions to run Maxent or  
325 similar species distribution models for their species of interest prior to the expedition in  
326 complement to expert knowledge on species distributions. Given sufficient input data,  
327 particular faith may be given to the model results for endemic species amongst the range of  
328 relevant crop wild relatives. The resulting maps should be combined with the knowledge of  
329 local authorities to identify potential new populations of these species. In addition, excluding  
330 occurrence data from the perceived range margins of the species may result in important  
331 information about local adaptation to distinct climatic conditions.

332 To increase access to the methodology, avoiding the necessity of installing and  
333 operating the software, a web-based version of Maxent, including the worldclim.org climatic  
334 dataset and standard model settings, would greatly facilitate the application of species  
335 distribution modelling in the preparation phase of collecting missions, and would be  
336 particularly useful for plant genetic resource conservation efforts with limited resources.

337

338 **ACKNOWLEDGMENTS**

339 This study was part of the Fundamental Research Programme on Sustainable Agriculture  
340 (KB-12-005.03-003) funded by the Dutch Ministry of Economic Affairs. M.C. was  
341 additionally funded by the Open Programme of the Netherlands Organisation of Scientific  
342 Research (NWO). Data and methods were contributed by the Adapting Agriculture to Climate  
343 Change Project, which is managed by the Global Crop Diversity Trust with the Millennium  
344 Seed Bank Partnership of the Royal Botanic Gardens, Kew and supported by the Government  
345 of Norway. For further information see the project website: <http://www.cwrdiversity.org/>

346 We acknowledge the contribution of the breeding companies Agrisemen BV, Enza  
347 Zaden Research and Development BV, Monsanto Holland BV, Nunhems Netherlands BV  
348 and Rijk Zwaan Breeding BV to the *Lactuca* expedition in 2013.



349 **REFERENCES**

350 Afonin A and Greene SL (1999) Germplasm collecting using modern geographic information  
351 technologies: directions explored by the N.I. Vavilov Institute of Plant Industry. In: Greene  
352 SL and Guarino L (eds) *Linking Genetic Resources and Geography: Emerging Strategies for*  
353 *Conserving and Using Crop Biodiversity*. CSSA Special Publication 27. American Society of  
354 Agronomy, Madison, Wisconsin. pp.75–85.

355 Alexander JM (2013) Evolution under changing climates: climatic niche stasis despite rapid  
356 evolution in a non-native plant. *Proceedings of the Royal Society B* 280: 20131446.

357 Araújo MB and Guisan A (2006) Five (or so) challenges for species distribution models.  
358 *Journal of Biogeography* 33: 1677-1688.

359 Araújo MB and New M (2007) Ensemble forecasting of species distributions. *Trends in*  
360 *Ecology and Evolution* 22: 42-47.

361 D'Andrea L, Broennimann O, Kozłowski G, Guisan A, Morin X, Keller-Senften J and Felber  
362 F (2009) Climate change, anthropogenic disturbance and the northward range expansion of  
363 *Lactuca serriola* (Asteraceae). *Journal of Biogeography* 36: 1573-1587.

364 Elith J, Graham CH, Anderson RP, Dudík M, Ferrier S et al. (2006) Novel methods improve  
365 prediction of species' distributions from occurrence data. *Ecography* 29: 129-151.

366 Elith J and Leathwick JR (2009) Species distribution models: ecological explanation and  
367 prediction across space and time. *Annual Review of Ecology, Evolution and Systematics* 40:  
368 677-697.

369 Greene SL, Hart TC and Afonin A (1999) Using geographic information to acquire wild crop  
370 germplasm for *ex situ* collections: I. map development and field use. *Crop Science* 39: 836–  
371 842.

372 Grime JP (1977) Evidence for the existence of three primary strategies in plants and its  
373 relevance to ecological and evolutionary theory. *The American Naturalist* 111: 1169-1194.

374 GRIN (2014) <http://www.ars-grin.gov/cgi-bin/npgs/html/taxon.pl?21359>

375 Guarino L (1995) Mapping the ecogeographic distribution of biodiversity. In: Guarino L,  
376 Ramanatha Rao V and Reid R (eds) *Collecting Plant Genetic Diversity Technical Guidelines*.  
377 CAB International, Wallingford, UK. pp.287–327.

378 Guisan A and Thuiller W (2005) Predicting species distributions: offering more than simple  
379 habitat models. *Ecology Letters* 8: 993-1009.

380 Hijmans RJ, Guarino L and Mathur P (2012) DIVA-GIS vs 7.5. A geographic information  
381 system for the analysis of species distribution data. Software and manual available at  
382 [www.diva-gis.org](http://www.diva-gis.org).

383 Hirzel AH, Hausser J, Chessel D and Perrin N (2002) Ecological-niche factor analysis: How  
384 to compute habitat-suitability map without absence data. *Ecology* 83: 2027-2036.

385 Jarvis A, Williams K, Williams D, Guarino L, Caballero PJ and Mottram G (2005) Use of  
386 GIS for optimizing a collecting mission for a rare wild pepper (*Capsicum flexuosum* Sendtn.)  
387 in Paraguay. *Genetic Resources and Crop Evolution* 52: 671–682.

388 Lebeda A, Dolezalová I, Feráková V and Astley D (2004a) Geographical distribution of wild  
389 *Lactuca* species (Asteraceae, Lactuceae). *The Botanical Review* 70(3): 328-356.

390 Lebeda A, Doležalová I, Astley D (2004b) Representation of wild *Lactuca* spp. (Asteraceae,  
391 Lactuceae) in world genebank collections. *Genetic Resources and Crop Evolution* 51: 167-  
392 174.

393 Lebeda A, Doležalová I, Křístková E, Kitner M, Petrželová I, Mieslerová B and Novotná A  
394 (2009) Wild *Lactuca* germplasm for lettuce breeding: current status, gaps and challenges.  
395 *Euphytica* 170:15–34.

396 Maxted N, Dulloo E, Ford-Lloyd BV, Iriondo JM, Jarvis A (2008) Gap analysis: a tool for  
397 complementary genetic conservation assessment. *Diversity and Distributions* 14:1018–1030.

398 Parra-Quijano M, Iriondo JM, Torres E (2012) Improving representativeness of genebank  
399 collections through species distribution models, gap analysis and ecogeographical maps.  
400 *Biodiversity and Conservation* 21: 79-96.

401 Phillips SJ, Anderson RP and Schapire RE (2006) Maximum entropy modeling of species  
402 geographic distributions. *Ecological Modelling* 190: 231-259.

403 Ramírez-Villegas J, Khoury C, Jarvis A, Debouck DG and Guarino L (2010) A gap analysis  
404 methodology for collecting crop genebanks: a case study with *Phaseolus* beans. *Plos ONE*  
405 5(10): e13497.

406 Stockwell D and Peters D (1999) The GARP modelling system: problems and solutions to  
407 automated spatial prediction. *International Journal of Geographical Information Science* 13:  
408 143-158.

409 Thuiller W, Lafourcade B, Engler R and Araújo MB (2009) BIOMOD—a platform for  
410 ensemble forecasting of species distributions. *Ecography* 32: 369-373.

411 Treuren R van, Coquin P and Lohwasser U (2012) Genetic resources collections of leafy  
412 vegetables (lettuce, spinach, chicory, artichoke, asparagus, lamb's lettuce, rhubarb and rocket  
413 salad): composition and gaps. *Genetic Resources and Crop Evolution* 59: 981-997.

414 VanDerWal J, Shoo LP, Graham C and Williams SE (2009) Selecting pseudo-absence data  
415 for presence-only distribution modeling: how far should you stray from what you know?  
416 *Ecological Modelling* 220: 589-594.

417 Zohary D (1991) The wild genetic resources of cultivated lettuce (*Lactuca sativa* L.).  
418 *Euphytica* 53: 31-35.

419

420 **TABLES**

421

422 **Table 1.** The species of the primary (I), secondary (II) and tertiary (III) gene pool of *Lactuca sativa* for  
 423 which we found occurrence data in herbaria and genebanks, respectively.

424

species	gene pool	herbarium samples	genebank accessions
<i>L. aculeata</i>	I	3	4
<i>L. altaica</i>	I	3	2
<i>L. azerbaijanica</i>	I	0	0
<i>L. dregeana</i>	I	4	0
<i>L. georgica</i>	I	17	1
<i>L. scarioloides</i>	I	1	0
<i>L. serriola</i>	I	23520	1177
<i>L. saligna</i>	II	1451	102
<i>L. virosa</i>	II	3318	102
<i>L. acanthifolia</i>	III	34	0
<i>L. aurea</i>	III	0	0
<i>L. longidentata</i>	III	0	0
<i>L. orientalis</i>	III	141	0
<i>L. quercina</i>	III	106	6
<i>L. sibirica</i>	III	854	2
<i>L. taraxacifolia</i>	III	2	0
<i>L. tatarica</i>	III	861	12
<i>L. viminea</i>	III	728	12
<i>L. watsoniana</i>	III	0	0

425

426

427 **Table 2.** The 19 bioclimatic variables used as input for the model, downloaded at a scale of 2.5 arcmin  
 428 from [www.worldclim.org](http://www.worldclim.org)  
 429

BIO1	Annual Mean Temperature
BIO2	Mean Diurnal Range (Mean of monthly (max temp - min temp))
BIO3	Isothermality (BIO2/BIO7) (* 100)
BIO4	Temperature Seasonality (standard deviation *100)
BIO5	Max Temperature of Warmest Month
BIO6	Min Temperature of Coldest Month
BIO7	Temperature Annual Range (BIO5-BIO6)
BIO8	Mean Temperature of Wettest Quarter
BIO9	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

430

431

432 **FIGURE LEGENDS**

433

434 **Figure 1.** The projected probabilities of occurrence for *L. saligna* in the Greek region when including  
435 and excluding the occurrence data of this region, on a scatter plot showing the changes in projected  
436 probabilities for the Greek region when omitting the occurrence data.

437

438 **Figure 2.** The projected probabilities of occurrence for *L. saligna* in the Israeli region when including  
439 and excluding the occurrence data in this region on a scatter plot showing the changes in projected  
440 probabilities for the Israeli region when omitting the occurrence data.

441

442 **Figure 3.** The projected probabilities of occurrence for *L. saligna* in the Eurasian region when  
443 including and excluding the occurrence data in this region on a scatter plot showing the changes in  
444 projected probabilities for the Eurasian region when omitting the occurrence data.

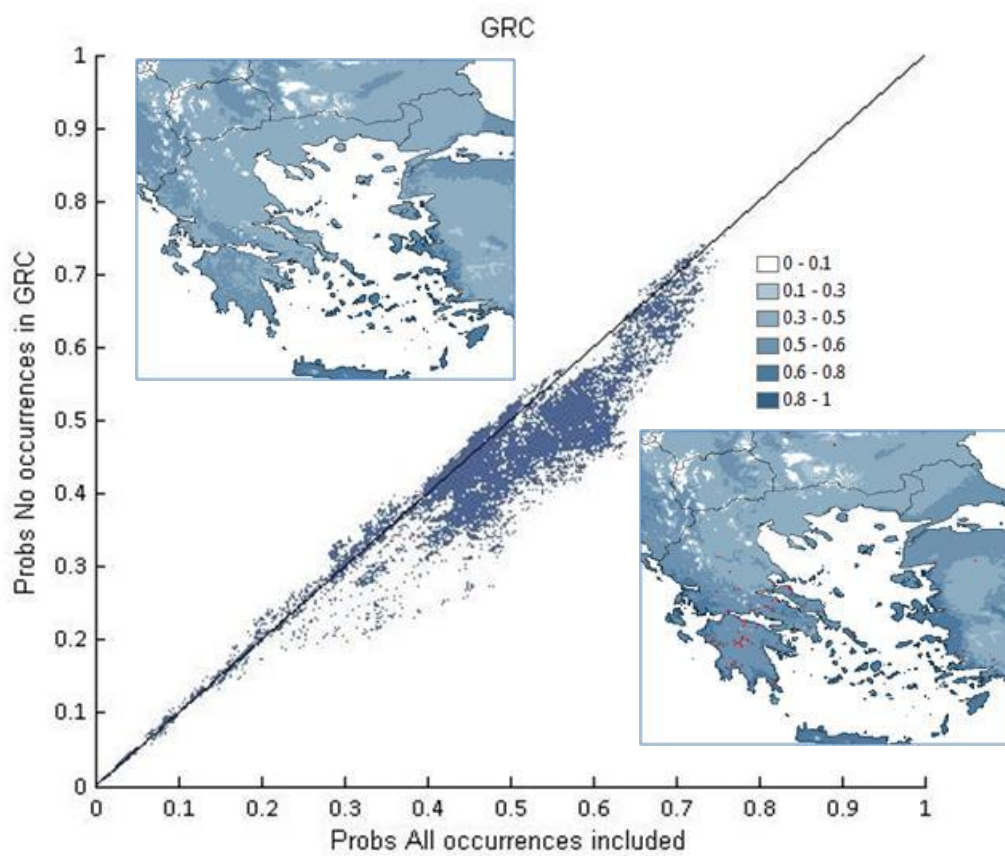
445

446 **Figure 4.** The number of locations in the region under study where *L. georgica* and *L. serriola* was  
447 present or absent plotted against the estimated probabilities of occurrence at these locations, based on  
448 the model including all occurrence data (*L. georgica* and *L. serriola*) and only occurrence data in the  
449 region under study (*L. serriola* CAU).

450

451 **FIGURES**

452 Figure 1

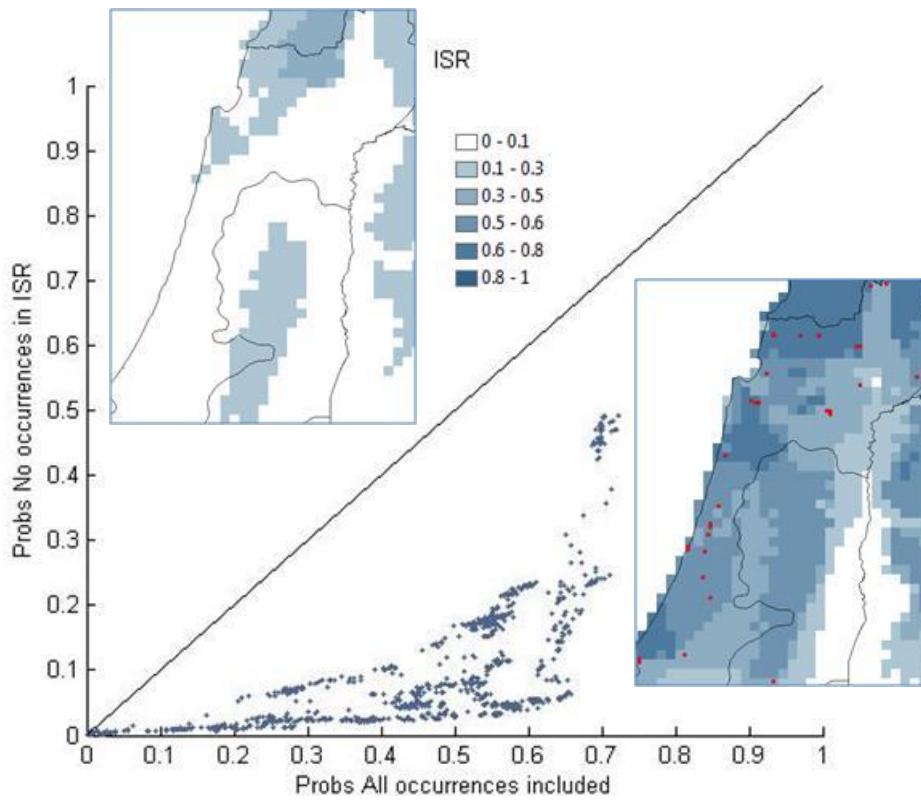


453

454



455 Figure 2

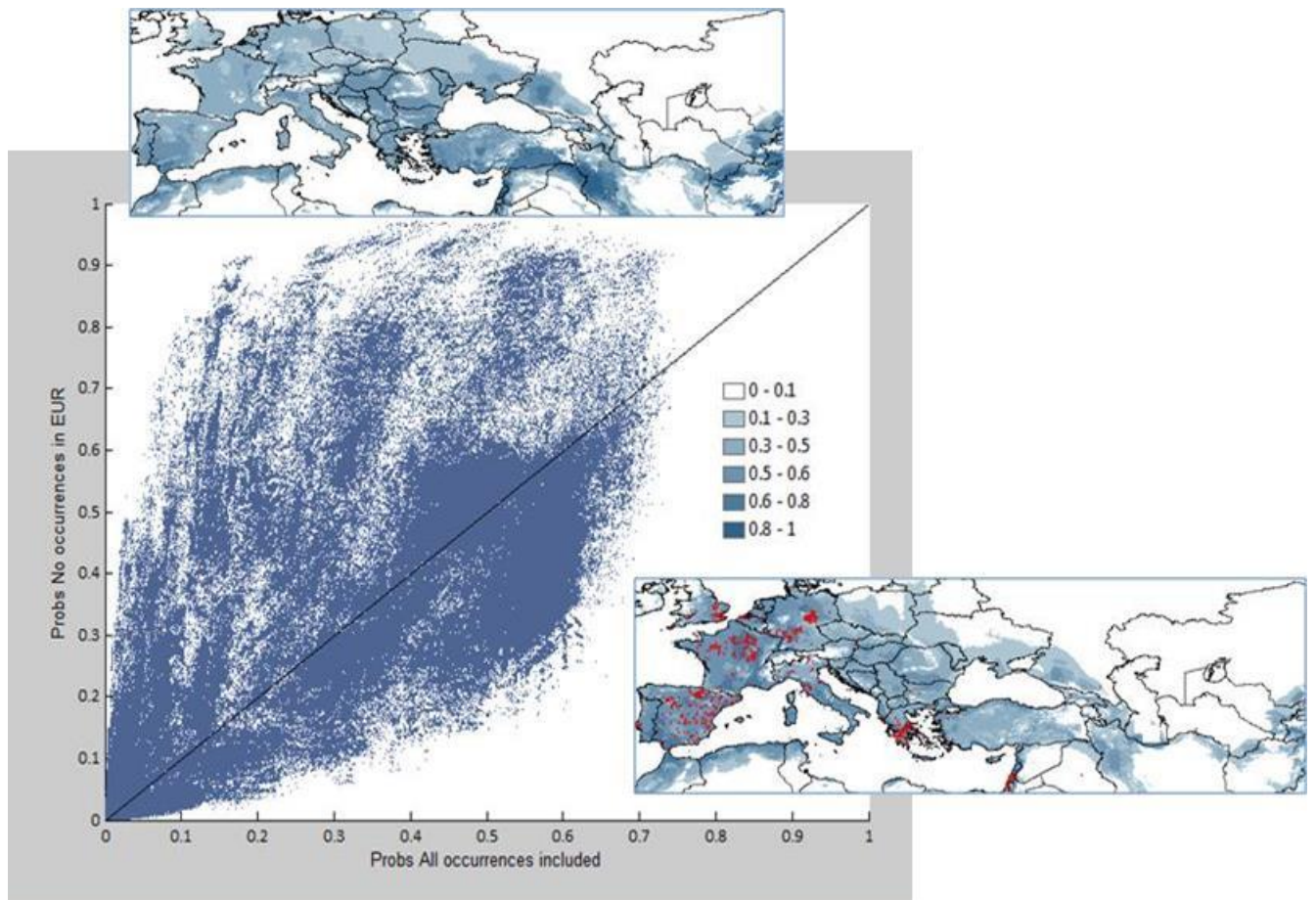


456

457

458

459 Figure 3



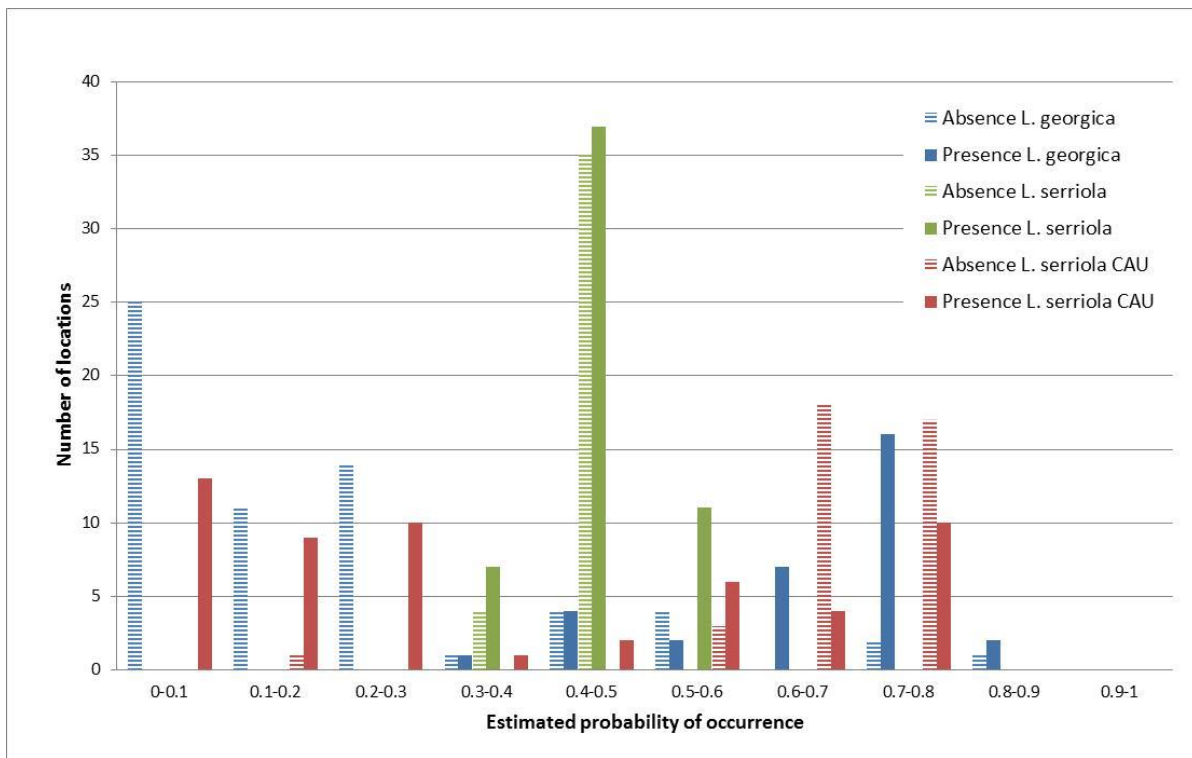
460

461

462

463

464 Figure 4



465

466

467

468

469 **SUPPLEMENT**

470 Table S1. Overview of all the herbaria and genebanks from which *Lactuca* samples were retrieved.

Academy of Natural Sciences Herbarium  
Australia's Virtual Herbarium  
Botanical Society of the British Isles Herbaria  
California Academy of Sciences  
Centro de Referência em Informação Ambiental  
Consortium of Pacific Northwest Herbaria  
Denver Botanic Gardens Herbarium  
Eurisco  
Field Museum  
Florida State University Herbarium  
Global Biodiversity Information Facility  
Harvard University Herbarium  
Instituto Superior de Agronomia  
International Lactuca Database  
Jardim Botânico do Rio de Janeiro  
Manchester University Herbarium  
Museum National d'Histoire Naturelle Herbarium  
Nationaal Herbarium Nederland  
New York Botanical Garden Herbarium  
Real Jardín Botánico de Madrid  
Royal Botanic Gardens Edinburgh  
Royal Botanic Gardens, Kew  
Smithsonian Institution, National Herbarium  
United States Department of Agriculture, National Plant Germplasm System  
Universidad del Valle Herbarium  
Universidade Lisboa Museu Nacional de História Natural e da Ciência  
University of California and Jepson Herbaria  
University of California, Riverside Herbarium  
University of Coimbra Herbarium  
V. L. Komarov Botanical Institute  
Wageningen University Herbarium  
West Virginia University Herbarium  
World Vegetable Center

471

472 Table S2. The species in the *Lactuca* genepool (Table 1) and their synonyms (Van Treuren *et al.*,  
 473 2012). Author names were included only when they distinguish between different synonyms in the  
 474 consulted databases (Table S1).

accepted names	synonyms
<i>L. acanthifolia</i>	<i>L. amorgina</i> <i>L. eburnea</i>
<i>L. aculeata</i>	
<i>L. altaica</i>	
<i>L. aurea</i>	
<i>L. azerbaijanica</i>	
<i>L. dregeana</i>	<i>L. virosa</i> Thunb.
<i>L. georgica</i>	
<i>L. longidentata</i>	<i>Scariola longidentata</i>
<i>L. orientalis</i>	
<i>L. quercina</i>	<i>L. altissima</i> <i>L. armena</i> <i>L. chaixii</i> <i>L. sagittata</i> <i>L. stricta</i> <i>L. vialea</i> <i>L. wilhelmsiana</i> <i>L. cracoviensis</i> <i>L. cyanea</i> <i>L. decorticata</i>
<i>L. saligna</i>	<i>L. adulteriana</i> <i>L. angustifolia</i> <i>L. caucasica</i> <i>L. cracoviensis</i> <i>L. cyanea</i> <i>L. salicifolia</i> <i>L. spiciformis</i> <i>L. tommasiniana</i> <i>L. virgata</i> <i>L. virosa</i> Habl. <i>L. wallrothii</i>
<i>L. scarioloides</i>	<i>L. kotschyana</i>
<i>L. serriola</i>	<i>L. albicaulis</i> <i>L. augustana</i> <i>L. coriacea</i> <i>L. dubia</i> <i>L. latifolia</i> <i>L. plicata</i> <i>L. scariola</i> <i>L. sylvestris</i>

	<i>L. tephrocarpa</i> <i>L. virosa</i> Luce
<i>L. sibirica</i>	
<i>L. taraxacifolia</i>	<i>L. alaica</i> <i>L. kotschy</i> <i>L. pentaphylla</i>
<i>L. tatarica</i>	<i>L. clarkei</i> <i>L. multipes</i> <i>L. oblongifolia</i> <i>L. pulchella</i> (Pursh) DC. <i>L. pulchella</i> DC.
<i>L. viminea</i>	<i>L. alpestris</i> <i>L. chondrilliflora</i> <i>L. decorticata</i> <i>L. numidica</i> <i>L. ramosissima</i>
<i>L. virosa</i>	<i>L. agrestis</i> <i>L. ambigua</i> <i>L. cornigera</i> <i>L. flavida</i> <i>L. lactucarii</i> <i>L. livida</i> <i>L. serratifolia</i> <i>L. sinuata</i> <i>L. virosa</i> L.
<i>L. watsoniana</i>	

475

476

477 Table S3. The Maxent model statistics for each of the investigated species of the lettuce genepool.  
 478 ATAUC: the 10-fold average test AUC (area under the curve), STAUC: the standard deviation of the  
 479 test AUC of the 10 different folds, ASD15: the percentage of the potential distribution coverage with  
 480 standard deviation above 0.15. For the gray colored species, the total number of samples is smaller  
 481 than 10 (Table 1), meaning that the number of folds is equal to the number of samples.

Taxon	ATAUC	STAUC	ASD15	ValidModel
<i>Lactuca acanthifolia</i>	0.9984	0.0006	0	yes
<i>Lactuca aculeata</i>	0.9599	0.0667	0	yes
<i>Lactuca altaica</i>	0.8728	0.0571	44.195	no
<i>Lactuca dregeana</i>	0.9896	0.0061	4.433	yes
<i>Lactuca georgica</i>	0.9973	0.0043	0	yes
<i>Lactuca orientalis</i>	0.9736	0.0240	0.610	yes
<i>Lactuca quercina</i>	0.9810	0.0341	0.055	yes
<i>Lactuca saligna</i>	0.9183	0.0044	0	yes
<i>Lactuca saligna_EUR</i>	0.9422	0.0048	0.037	yes
<i>Lactuca saligna_GRC</i>	0.9208	0.0057	0	yes
<i>Lactuca saligna_ISR</i>	0.9191	0.0051	0	yes
<i>Lactuca scarioloides</i>	NA	NA	NA	no
<i>Lactuca serriola</i>	0.6490	0.0046	0	no
<i>Lactuca serriola_TC</i>	0.9895	0.0060	1.450	yes
<i>Lactuca sibirica</i>	0.9596	0.0038	0	yes
<i>Lactuca taraxacifolia</i>	0.5000	0	NA	no
<i>Lactuca tatarica</i>	0.9169	0.0061	0.126	yes
<i>Lactuca viminea</i>	0.9637	0.0043	0.003	yes
<i>Lactuca virosa</i>	0.8806	0.0042	0	yes

482

483 Figure S1a. The global region of analysis (ALL) and all *L. saligna* occurrences.

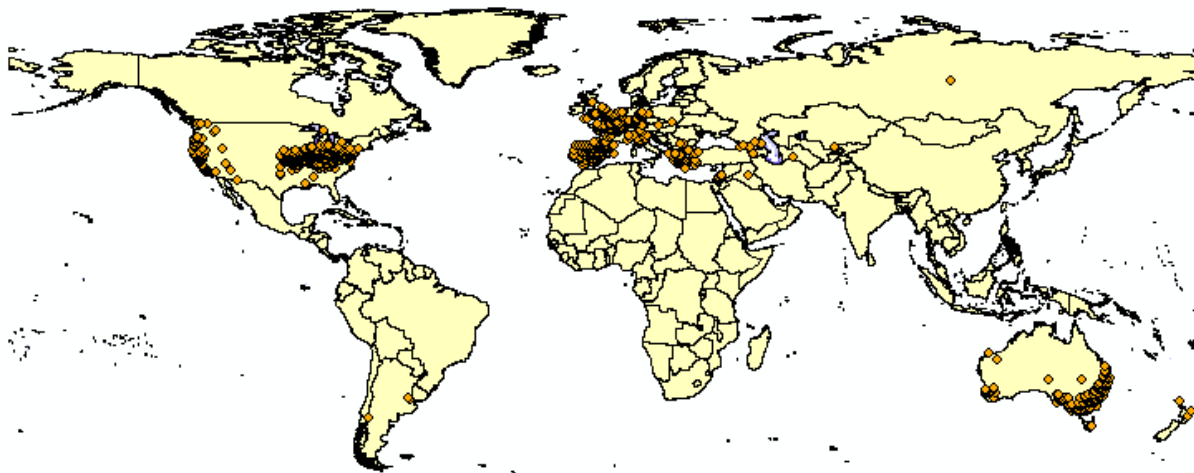
484 Figure S1b. The Eurasian region of analysis (EUR) and its *L. saligna* occurrences.

485 Figure S1c. The Greek region of analysis (GRC) and its *L. saligna* occurrences.

486 Figure S1d. The Israeli region of analysis (ISR) and its *L. saligna* occurrences.

487

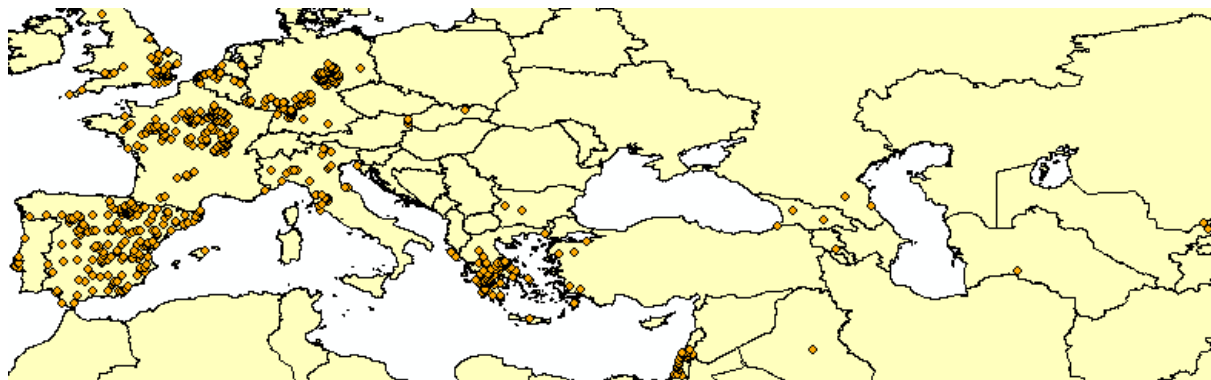
488 Figure S1a



489

490

491 Figure S1b



492

493



494

495 Figure S1c

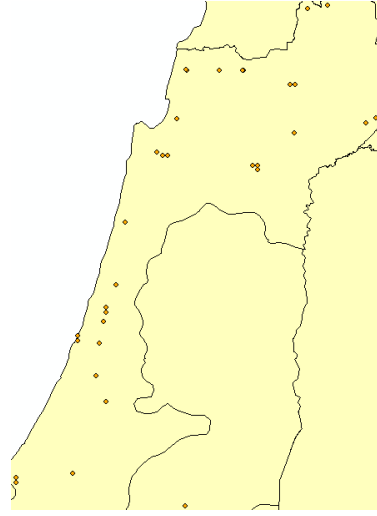


496

497

498

Figure S1d



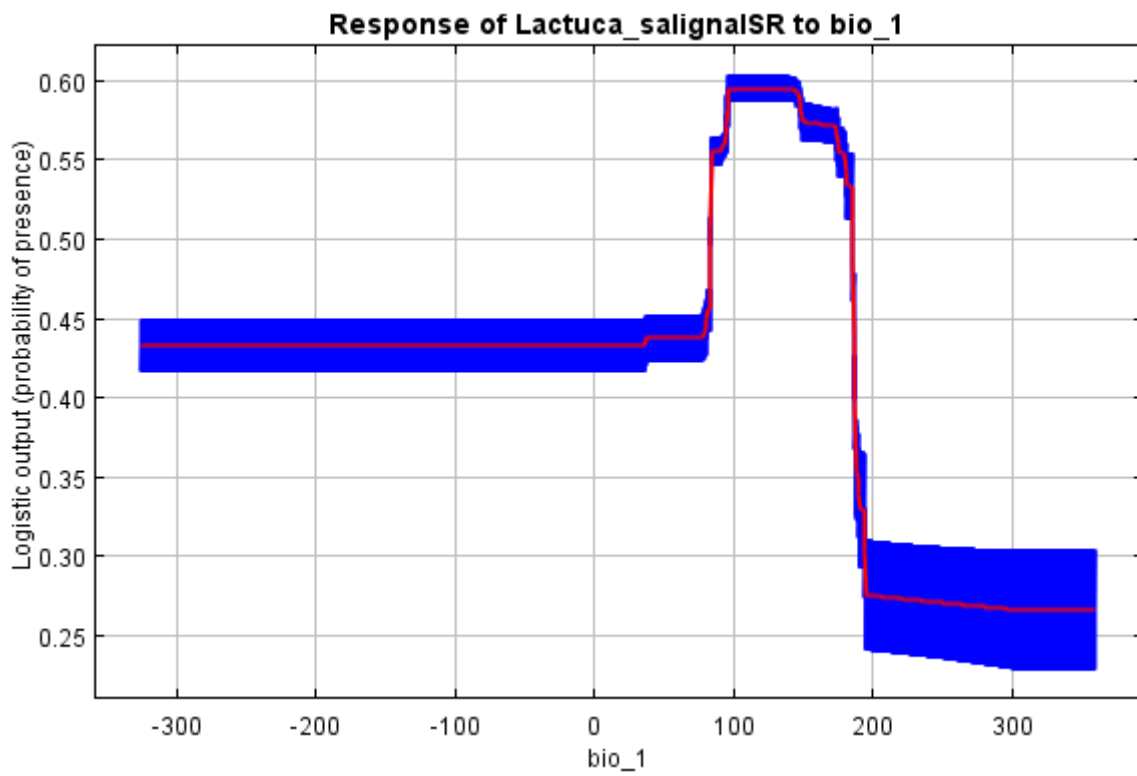
499 Figure S2a. The Maxent response curves to BIO\_1 for the model based on all occurrences  
500 (model *Lactuca\_saligna*) and the model for which the Israeli occurrences were omitted (model  
501 *Lactuca\_salignaISR*). The curves show how the logistic prediction changes as the mean  
502 annual temperature is varied, keeping all other environmental variables at their average  
503 sample value.

504 Figure S2b. The Maxent response curves to BIO\_5 for the model based on all occurrences  
505 (model *Lactuca\_saligna*) and the model for which the Israeli occurrences were omitted (model  
506 *Lactuca\_salignaISR*). The curves show how the logistic prediction changes as the maximum  
507 temperature of the warmest month is varied, keeping all other environmental variables at their  
508 average sample value.

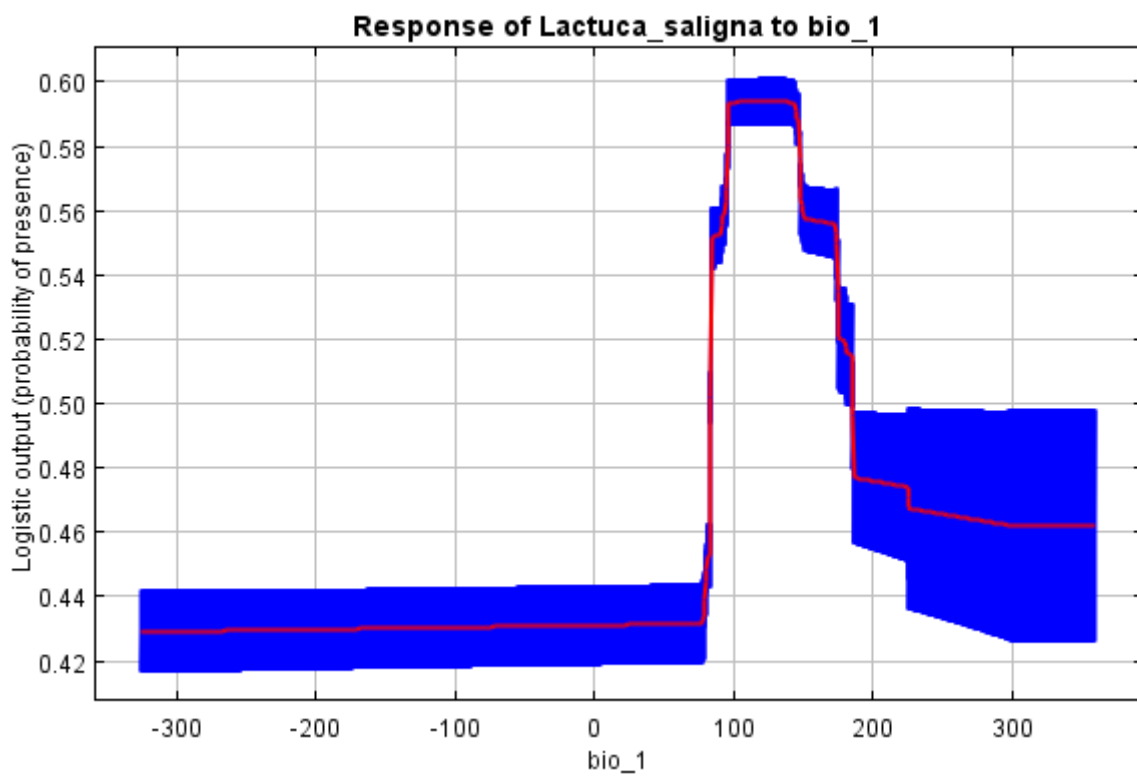
509 Figure S2c. The Maxent response curves to BIO\_15 for the model based on all occurrences  
510 (model *Lactuca\_saligna*) and the model for which the Israeli occurrences were omitted (model  
511 *Lactuca\_salignaISR*). The curves show how the logistic prediction changes as the variation in  
512 the precipitation over the seasons is varied, keeping all other environmental variables at their  
513 average sample value.

514 Figure S2d. The Maxent response curves to BIO\_18 for the model based on all occurrences  
515 (model *Lactuca\_saligna*) and the model for which the Israeli occurrences were omitted (model  
516 *Lactuca\_salignaISR*). The curves show how the logistic prediction changes as the total  
517 precipitation of the warmest quarter is varied, keeping all other environmental variables at  
518 their average sample value.

519 Figure S2a



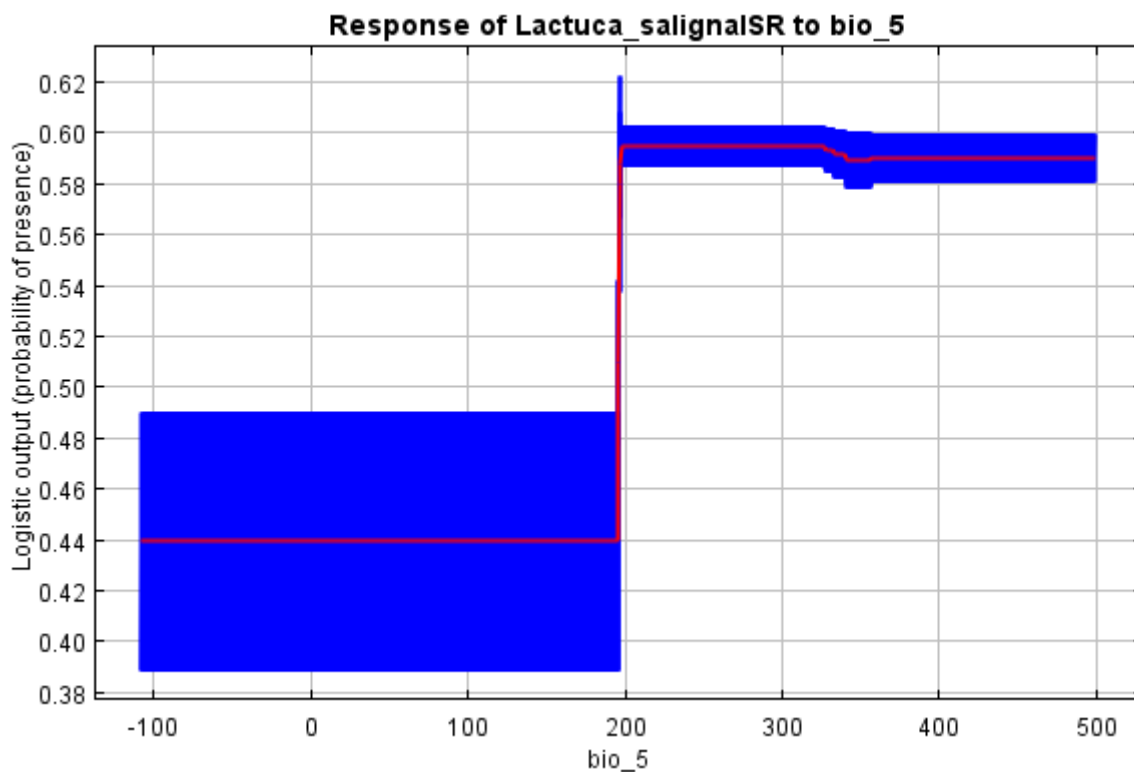
520



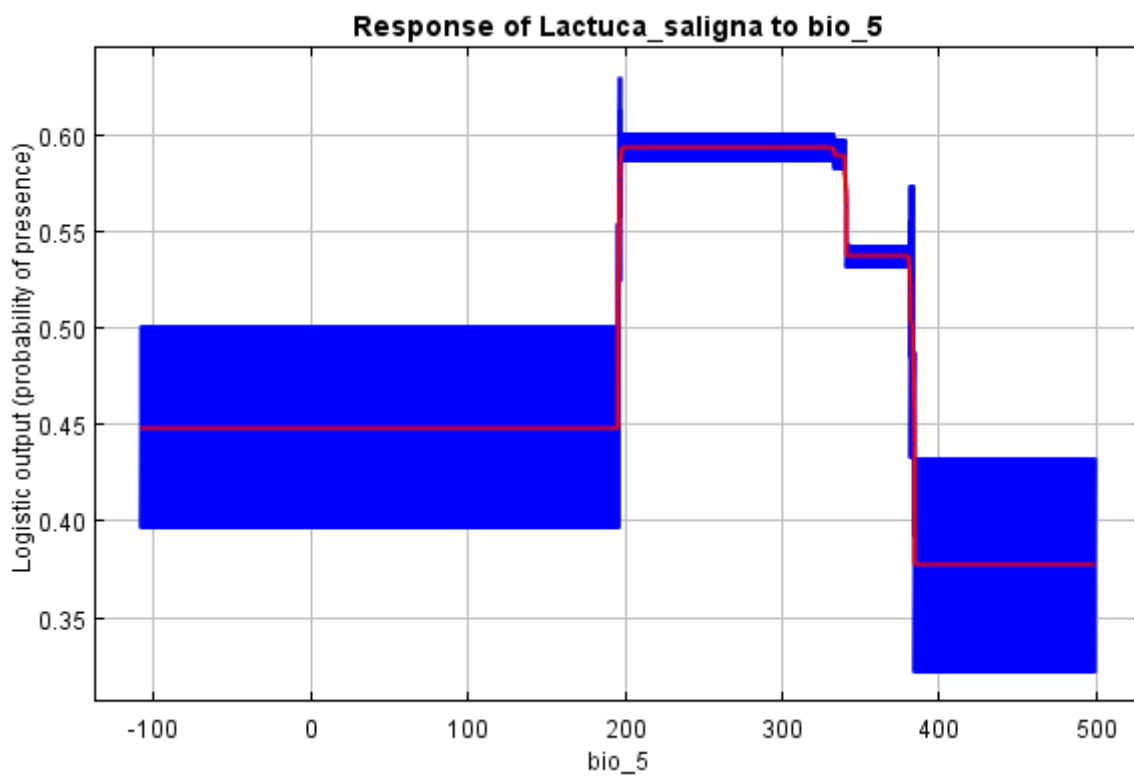
521

522

523 Figure S2b



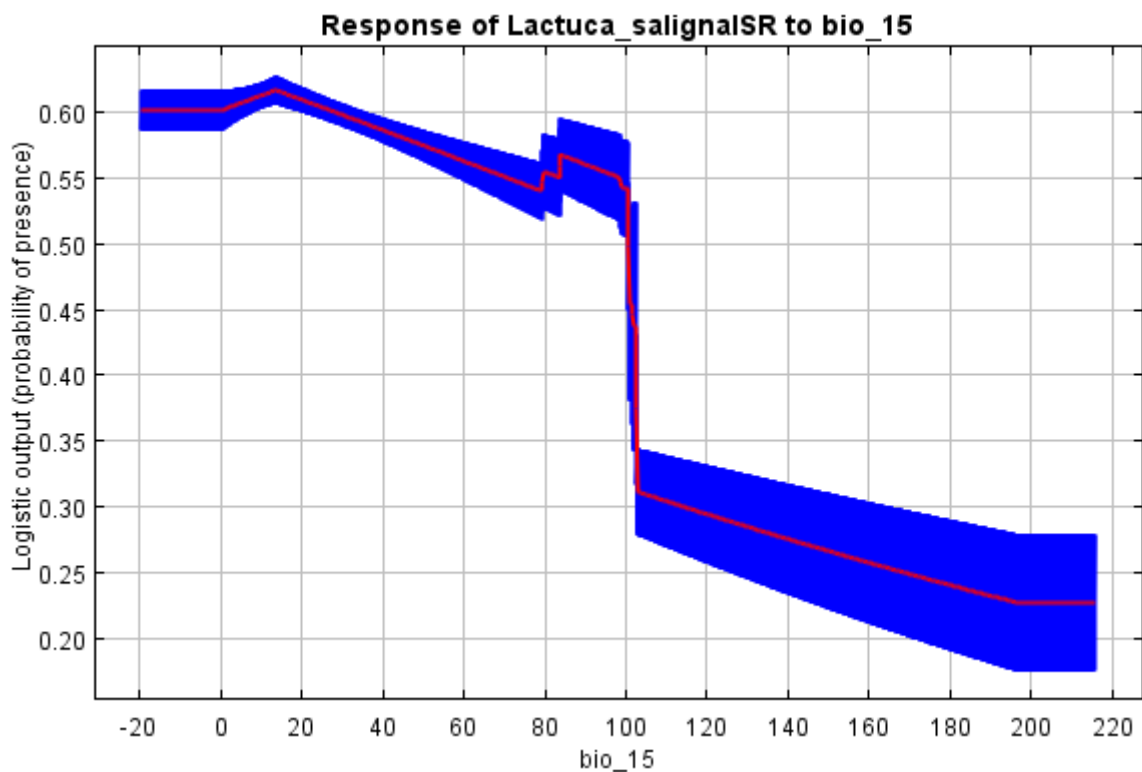
524



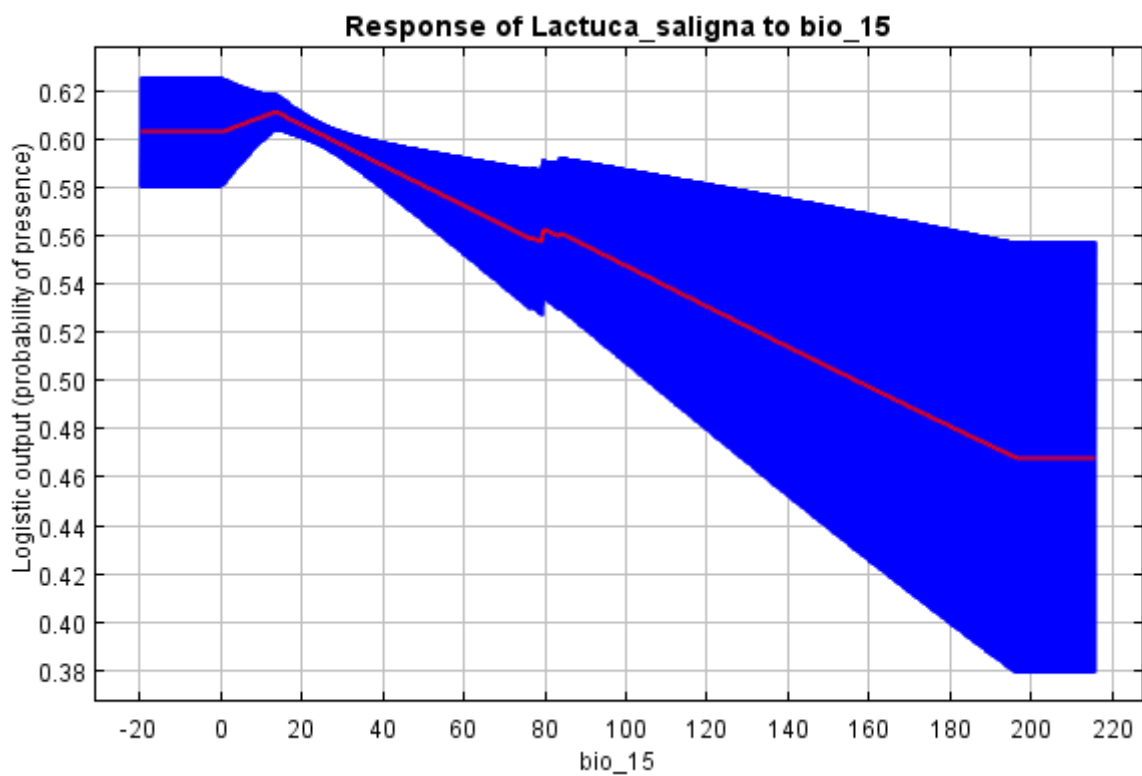
525

526

527 Figure S2c



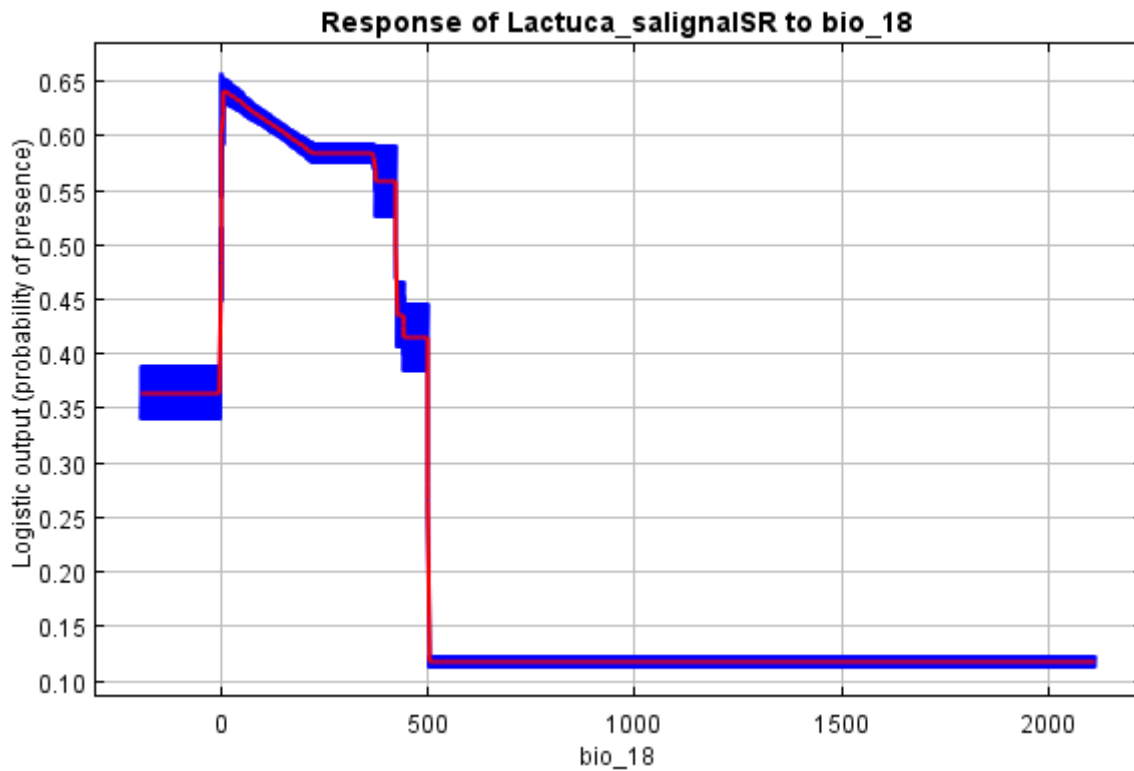
528



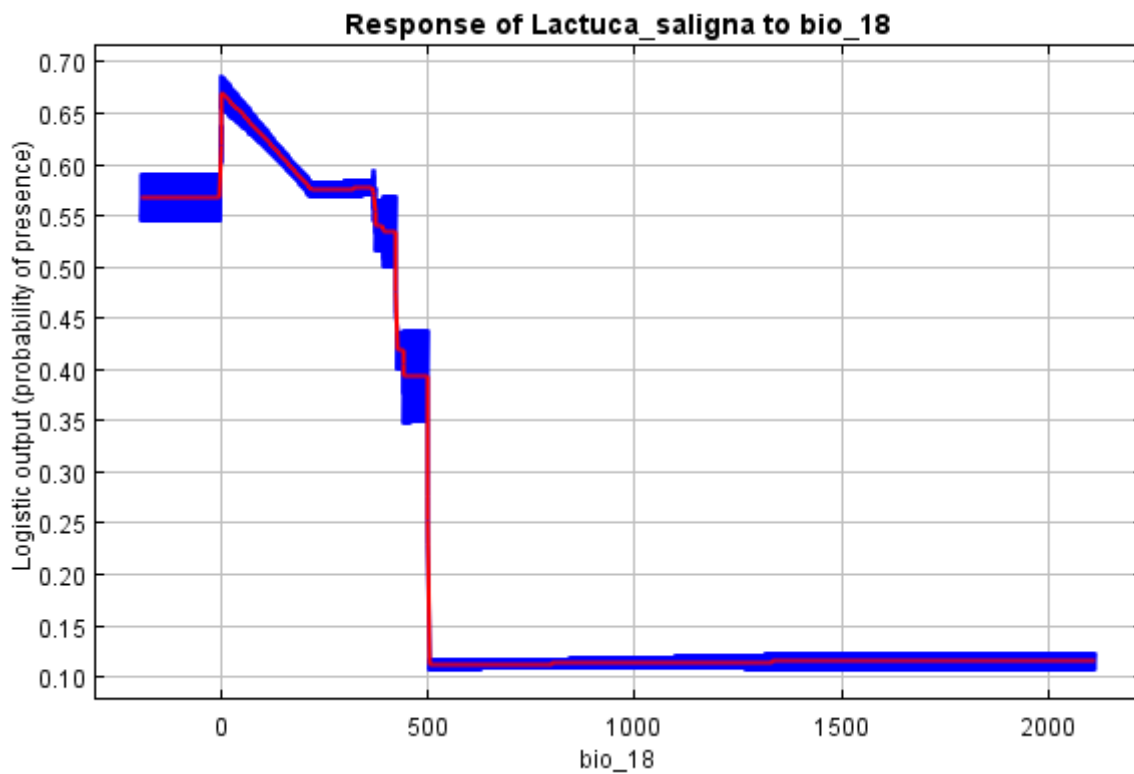
529

530

531 Figure S2d



532



533

534