# How to assess species distribution model accuracy: using internal-aspatial or external-spatial methods?

Chunrong Mi [1, 2, 3] , Falk Huettmann [4] , Yumin Guo [Corresp. 3]

[1] Chinese Academy of Science, Key Laboratory of Animal Ecology and Conservation Biology, Institute of Zoology, Beijing, China

[2] University of Chinese Academy of Science, Beijing, Beijing, China

[3] Beijing Forestry University, College of Nature Conservation, Beijing, China

[4] University of Alaska Fairbanks (UAF), EWHALE Lab- Department of Biology and Wildlife, Institute of Arctic Biology, Fairbanks, Alaska, USA

Corresponding Author: Yumin Guo
Email address: guoyumin@bjfu.edu.cn

Species distribution models (SDMs) have become an increasingly important tool in ecology, biogeography, evolution and, more recently, in conservation management, landscape planning and climate change research. The assessment of their predictive accuracy is one fundamental issue in the development and application of SDMs. Accuracy assessments for models should have a close connection to the intended use of the model. However, we found that the common evaluation method (we named internal-aspatial) usually ignored how the spatial prediction map actually looks like, and achieves for the real-world species distribution and for application. Therefore, in this research we proposed a spatial method to evaluate model performance by assessing how the prediction maps look like (we named external-spatial). We took Hooded Crane (*Grus monacha*) as a case, in this research, to compare these two methods (internal-aspatial and external-spatial) performance. Both of the two methods were expressed with three commonly used SDM evaluation criteria (AUC, Kappa and TSS). In addition, model accuracy was also assessed via evaluating the prediction maps with knowledge of the study species and alternative occurrence data assistance. We used two popular data mining algorithms (Random Forest and TreeNet) and ran 8 experiments using 1, 3, 5, 8, 11, 21, 29 and 78 predictors, allowing to develop overall 16 models for this assessment. Results indicated that AUC had a significant linear relationship with Kappa and TSS. Both of interal-aspatial and external-spatial methods could get higher AUC values and they were close. This indicated that internal-aspatial model assessments can serve as powerful assessment-aspatiual metrics without the need of secondary data even! However, internal-aspatial, external-spatial, prediction map evaluation and alternative occurrence data could not distinguish well models with different sets of predictors. This is the first time the concept of spatial assessment criteria is expressed and assessed. Overall, we hope to see more study on meaningful spatial criteria and proposed more and better methods to evaluate SDMs and

distribution map in the future.

1    How to assess species distribution model accuracy: using internal-

2    aspatial or external-spatial methods?

3    Chunrong Mi[1,2,3], Falk Huettmann[4], Yumin Guo[3]

4    [1]Key Laboratory of Animal Ecology and Conservation Biology, Institute of Zoology, Chinese

5    Academy of Sciences, Beijing, China

6    [2]University of Chinese Academy of Sciences, Beijing, China

7    [3]College of Nature Conservation, Beijing Forestry University, Beijing 100083, China

8    [4]EWHALE Lab- Department of Biology and Wildlife, Institute of Arctic Biology, University of

9    Alaska Fairbanks (UAF), 419 Irving I, Fairbanks Alaska AK 99775, USA

10

11    Corresponding author:

12    Yumin Guo

13    College of Nature Conservation, Beijing Forestry University, P.O. Box 159, Beijing 100083,

14    China

15    guoyumin@bjfu.edu.cn

16

17

18

19

20

21

22

23

24

25

26

## ABSTRACT

Species distribution models (SDMs) have become an increasingly important tool in ecology, biogeography, evolution and, more recently, in conservation management, landscape planning and climate change research. The assessment of their predictive accuracy is one fundamental issue in the development and application of SDMs. Accuracy assessments for models should have a close connection to the intended use of the model. However, we found that the common evaluation method (we named internal-aspatial) usually ignored how the spatial prediction map actually looks like, and achieves for the real-world species distribution and for application. Therefore, in this research we proposed a spatial method to evaluate model performance by assessing how the prediction maps look like (we named external-spatial). We took Hooded Crane (*Grus monacha*) as a case, in this research, to compare these two methods (internal-aspatial and external-spatial) performance. Both of the two methods were expressed with three commonly used SDM evaluation criteria (AUC, Kappa and TSS). In addition, model accuracy was also assessed via evaluating the prediction maps with knowledge of the study species and alternative occurrence data assistance. We used two popular data mining algorithms (Random Forest and TreeNet) and ran 8 experiments using 1, 3, 5, 8, 11, 21, 29 and 78 predictors, allowing to develop overall 16 models for this assessment. Results indicated that AUC had a significant linear relationship with Kappa and TSS. Both of interal-aspatial and external-spatial methods could get higher AUC values and they were close. This indicated that internal-aspatial model assessments can serve as powerful assessment-aspatiual metrics without the need of secondary data even! However, internal-aspatial, external-spatial, prediction map evaluation and alternative occurrence data could not distinguish well models with different sets of predictors. This is the first time the concept of spatial assessment criteria is expressed and assessed. Overall, we hope to see more study on meaningful spatial criteria and proposed more and better methods to evaluate SDMs and distribution map in the future.

***Keywords:*** Species distribution models (SDMs); internal-aspatial; external-apatial; AUC; Kappa;

53     TSS; Random Forest; TreeNet (Boosting); Hooded Crane (*Grus monacha*)

## INTRODUCTION

Species distribution models (SDMs) are well-established numerical spatial tools that combine field observations of species occurrence or abundance with environmental predictor variables (Guisan and Zimmermann, 2000; Elith and Leathwick, 2009; Drew et al., 2011). In recent years, predictive model of species distribution has become an increasingly important science-based management tool for policy (e.g. Drew et al., 2011; Mi et al., 2016). The trend of SDMs goes towards data mining of data heavy applications to address many wider and more holistic and interdisciplinary issues in ecology, biogeography, evolution and, more recently, in conservation management and climate change research, and usually gets done now on a global level but with a high resolution (Austin et al., 1990; Franklin, 1995; Peterson et al., 2002; Guisan and Thuiller, 2005; Wei et al., 2010). A wide variety of statistical and machine learning methods have been introduced, often in conjunction with geographic information systems (GIS), remote-sensing (Aspinall and Veitch, 1993; Franklin et al., 2000; Hegel et al., 2010).

Only slowly, a number of 'data mining' approaches for modeling data that contain non-linear and other complex and interacting dependencies have appeared now in the wildlife literature, too (Derrig and Francis, 2006; Cushman and Huettmann 2010; Mi et al 2014). Relatively new methods are based on data hungry networks, or ensembles and that can handle complex and even marginal data situations. Over 100 machine learning algorithms exist (Fernandez-Delgado et al. 2014). Some of them use statistical trees and include algorithms such as Random Forest (bagging; Breiman, 2001a), TreeNet (stochastic gradient boosting (Friedman, 2002)), and other methods (Araujo and New 2007; see Biomod2 package in R). Random Forest and TreeNet were used in this research, because model construction process was fast and convenient (Mi et al., 2014 and 2016), and also model prediction performed very well in similar investigations (Oppel et al., 2012, Mi et al., 2014), and they were also non-parametric. These two algorithms usually do not really require or make relevant *a priori* assumptions about the relationship between the response and predictor variables. It does not limit the number of predictor variables, and it is capable of uncovering the underlying structure of data that are non-additive, interacting or hierarchical in nature (Prasad et al., 2006; Hasti et al., 2009; Cushman and Huettmann, 2010; Drew et al., 2011).

While somewhat overlooked, the deeper and ecological assessment of predictive accuracy is

84    one fundamental issue in machine learning and specifically in the development of species
85    distribution models (Fielding and Bell, 1997; Pearce and Ferrier, 2000; Guisan and Thuiller,
86    2005; Allouche et al., 2006). A quantitative assessment of model performance assists in
87    determining the suitability of the model for specific applications (Vaughan and Ormerod, 2005;
88    Barry and Elith, 2006; Guisan et al., 2006). Model performance assessment can also provide a
89    basis for comparing alternative modelling techniques (Loiselle et al., 2003; Segurado and Araujo,
90    2004; Pearson et al., 2006) and it enables the user to investigate how different properties of the
91    data and/or the species affect the accuracy of predictive maps generated by the model (Kadmon
92    et al., 2003; Segurado and Araujo, 2004; Reese et al., 2005; Seoane et al., 2005).

93        To assess model performance, we found that most scholars used evaluation criteria (e.g. the
94    Area Under the ROC Curve (AUC), the Kappa Statistic (Kappa) and the True Skill
95    Statistic(TSS)) were created by the software itself (e.g. Anderson and Gonzalez Jr, 2011; Elith et
96    al., 2011). Typically, such tests are based on hold-out data, such as from boot strapping or
97    jackknifing applied to the (large) training data. This creates unmapped (aspatial) models and then
98    calculates metrics (Kappa, TSS). We named this: internal-aspatial method. However, this method
99    ignored how the spatial prediction map actually looks like, and achieves for the real-world
100   species distribution and for application. Actually, accuracy assessments for models should have a
101   close connection to the intended use of the model (Fielding, 2002). One major role of species
102   distribution models is to model complex ecology and to support an efficient conservation
103   management, such as conservation planning, design reserve networks that maintain biodiversity
104   (Guisan and Thuiller, 2005). Therefore, in this research we proposed a spatial method to evaluate
105   model performance by assessing how the prediction maps look like. We named this: external-
106   spatial method. This spatial metrics mean it mapped predictions first, and then used presence-
107   absence points to overlay the created prediction map and get the relative index of occurrence
108   (RIO), and then calculated accuracy metrics to obtain the estimate of accuracy, such as AUC,
109   Kappa and TSS. The aim of this research was to explore which evaluation method was better,
110   internal-aspatial or external-spatial? In addition, we also assess model accuracy using prediction
111   map with experts' knowledge of target species' distribution, and alternative occurrence data.

## MATERIALS AND METHODS

### Study species put to a test

The Hooded Crane is listed as a vulnerable (VU) species in the IUCN Red List. This species breeds in Eastern Russia and Northeastern China (Guo, 2005; Simonov and Dahmer, 2008; Mi et al., 2018). Its global population is estimated to be 11,160 individuals (Birdlife international, 2014) and the population size is declining (IUCN, 2012). In recent years, more than 10,500 (~ 94%) Hooded Cranes winter in Izumi, Japan (Birdlife international, 2014). This presents a risk and therefore, it is badly needed to find suitable places and methods to disperse the Hooded Crane from Izumi in order to diversify and reduce the population density there and to minimize local risks. Otherwise it can for instance lead to epidemic diseases of birds and their population crashes, such as Avian Influenza (Mi et al., 2018). Thus, here we tried to construct a winter distribution model for Hooded Cranes and to see where they would stay and for making a conservation plan and obtain more management methods.

### Species occurrence data

The Hooded Crane winter occurrence data was collected from our own fieldwork, also using previously published literature in East Asia (Fig. 1). In general, the data were initially provided by location name. To ensure the exact position for a valid geo-referencing, we then searched the location using a map with coordinates, Google Earth and also consulted experts for confirmation. Overall, we obtained 112 data points that were observed for this species during 1980-2013. This compiled data represents the best available geo-referenced data set for Hooded Crane wintering in China we know (Supplement S1). Initially, we considered that the data points maybe overly dense in some locations (oversampled or cluster sampling); thus, we created a concentric buffer around a data point with a 2-km radius in ArcGIS 10.1 (Toolboxes/System Toolboxes/Analysis Tools/Proximity/Buffer). However, in our data, we did not find any overlap of the 2-km scale. Therefore, we continued to use the all the data points as intended.

### Environmental layers

The environmental predictor variables we used to develop models in this study describe

139   climate, topography, terrain and human factors. We chose a set of 21 predictors (see Supplement

140   S2) to develop models as the ordinary baseline model, which was often used in our previous

141   research (e.g. Mi et al., 2017). Based on this step, we then added another 8 bio-climate layers

142   (Supplement S2) to construct models with overall 29 predictors (Mi et al., 2016; Han et al.,

143   2018). In addition, we tried to make models with the entire predictor set (78 predictors). The

144   Salford Predictive Modeler (SPM) suite we applied can generate a variable importance ranking

145   table from the obtained trees; here we chose the top 1, top 3, top 5, top 8 and top 11 predictor

146   variables. For that initial approach, we refer readers to Harrell et al. (1996), who promote for

147   multiple regressions that the variable (predictor) quantity should not exceed n/10 (n means

148   sample size, in our case n=112) for multivariable regression models). In all, we created 8

149   Random Forest models and 8 TreeNet models. All data layers were publicly available and had a

150   global-wide coverage (Supplement S2). We re-projected layers into WGS-1984 Mercator (in

151   meters) and merged them for a study area coverage in ArcGIS. Slope and aspect layers were

152   derived in ArcGIS from the DEM. We also calculated the Euclidean distance to road, railroad,

153   river, lake, coastline, settlement using the Euclidean distance tool in ArcGIS 10.1. Layers and

154   raw data can be obtained from the College of Nature Conservation, Beijing Forestry University

155   upon request to the autuors.

156   <span style="color:red">Put Figure 1 Here</span>

## Selection of model algorithms

157

158         In SPM, we choose Random Forest (hereafter RF) and TreeNet (hereafter TN) as our

159   species distribution models. We used them as a set of representative algorithms for the wider

160   machine learning (ML) family of methods. RF and TN are specific stand-alone software products

161   from Salford Systems Ltd that can outcompete R implementations (Herrick 2013), and each

162   performs one specific technique. Here we used them as representative ML methods because

163   when using these algorithms, model construction is fast and convenient, they offer a very high

164   degree of fault tolerance for messy and incomplete data (Friedman, 2001; Craig and Huettmann

165   2008, Mi et al., 2014; Jiao et al. 2016). For more details on Random Forest and TreeNet, we refer

166   readers to the user guide (https://www.salford-systems.com/products/spm/userguide).

## Model development

168    We created 10,000 random points across the study area using the freely available Geospatial
169    Modeling Environment (GME; Beyer, 2013; Booms et al., 2010) and compared these points to
170    the 112 bird locations. The ratio of 112 presence points versus 10,000 pseudo-absence points is
171    commonly used in the machine learning modeling literature (Booms et al., 2010, Mi et al., 2016)
172    and even more so when it comes to data mining (Hastie et al. 2009, Jiao et al. 2016) We
173    extracted information from environmental layers at bird location sites and random points using
174    GME.

175    We generally used the powerful default settings in SPM (e.g. Mi et al., 2014). Our
176    distribution models were constructed in SPM by using 'classification' and the balanced class
177    weights option to account for unequal sample sizes of presence and availability (pseudo-absence).
178    For the predictions, we created equally-spaced point lattice grids of 1,047,746 regularly spaced
179    points across our study area (approximately a 5×5 km spacing for the study area). We extracted
180    information from the environmental layers (Supplement S2) described above for each point, and
181    then used the model to predict (='score') birds occurrence as a relative index of occurrence at
182    each lattice point based on the extracted environmental data. For visualization, we imported the
183    dataset of spatially referenced predictions into GIS as a raster file, and interpolated for visual
184    purposes between the regular points using inverse distance weighting (IDW) to obtain a
185    smoothed predictive map, as it is commonly done (e.g. Kandel et al. 2015, Regmi et al. 2018).
186    We used that resulting prediction surface for our spatial assessment and comparison.

## Evaluation criteria

188    In this study, we obtained three metrics: AUC (the area under the ROC Receiver Operator
189    Curve), the Kappa Statistic (Kappa) and the True Skill Statistic (TSS). They were among the
190    most popular measures that are commonly used to assess the accuracy of prediction models
191    (Fielding and Bell, 1997; Pearce and Ferrier, 2000; Manel et al., 2001; Pearson et al., 2004;
192    Huettmann and Gottschalk, 2011; Liu et al., 2013). They are based on the confusion matrix
193    (Fielding and Bell, 1997; Pearce and Ferrier, 2000), and AUC is more generally applied whereas
194    Kappa and TSS offer specific advantages in terms of prevalence (Manel et al., 2001). None of
195    those metrics take the spatial distribution, arrangement or autocorrelation of the points into

196 account though (Betts et al. 2009). We obtained these three criteria using two methods (internal-
197 aspatial and external-spatial). When SPM creates the model, it offers internal-aspatial AUC value,
198 and a threshold table which could be transformed as a subsequent confusion matrix table, then
199 we calculated internal–aspatial Kappa and TSS in R 3.0.3 according to the formula in
200 Supplement S3. For the spatial metric, we extracted the relative index of occurrence (RIO) for all
201 of the presence and pseudo-absence points from each model prediction map with GME software.
202 Using the above RIO and the "SDMTools" package in R 3.0.3, we could obtain external-spatial
203 AUC, Kappa, and TSS. Kappa and TSS values can be shown with different thresholds (0-1,
204 interval 0.01), in this research, we used the maximum value (max-Kappa, max-TSS) as the final
205 metric.

## Prediction map assessment

207 Model prediction maps were evaluated by the 'true' distribution of Hooded Cranes in winter,
208 as we know them from our own field experience. We ranked the prediction map based on
209 following reasons:
210 (1) Closeness between the predicted distribution and the distribution we knew;
211 (2) Whether the predicted distribution reflects ecological and biology realities of Hooded Crane
212 (such as food, water availability etc);
213 (3) Whether some places were not predicted as part of the distribution area, and some places
214 should be not considered as the distribution area but they were predicted (such as a
215 settlement for instance)

## Alternative occurrence data assessment

217 Alternative data from other sources, such as other research, citizen observtions, speciemen,
218 or new field investiantion data were very important testing data for us to assess model accuracy
219 (Magness et al., 200; Huettmann and Gottschalk, 2011). In this study, we used the occurrence
220 data of Hooded Cranes from Global Biodiversity Information Facility (GBIF,
221 http://www.gbif.org/) as alternative data to assess our model performance. We applied all of the
222 90 locations, which were observed by people and recorded GPS location in winter time (October
223 to next February) from 1994 to 2013. Then we extracted the relative index of occurrence (RIO)
224 for each point from 8 Random Forest and 8 TreeNet spatial distribution map.

# RESULTS

## Static metric evaluation

AUCs and TSSs from both, internal-aspatial and external-spatial metrics were close among models with different number of variables in Random Forest, while the Kappa value diversified more between different models. The internal-aspatial AUCs were slightly greater than the related external-spatial metric; however, TSSs had a contrasting trend. For TreeNet models, internal-aspatial AUCs and TSSs were always larger than related external-spatial AUCs and TSSs (Fig. 2a and 2c). For the Kappa statics of Random Forest, it showed a somewhat contrasting result from RF3 to RF78 for internal and external metrics. The trend of internal metrics was increasing first and then decreasing; while external-spatial metrics kept increasing, except for RF21 and RF29 they were smaller than RF11 (Fig. 2b).

From the linear regression analysis (Table 1), we found that AUC had a significant positive relationship with TSS and Kappa ($P \leqslant 0.001$, $R^2 > 0.510$), both for internal-aspatial and external-spatial metrics, except for the internal-aspatial Kappa metric of the Random Forest model. Therefore, in the remaining analysis, we just used AUC to evaluate model accuracy.

<span style="color:red">Put Figure 2 Here</span>

<span style="color:red">Put Table 1 Here</span>

We used three-way ANOVA analysis between AUC and model algorithm (RF or TN), evaluation method (internal-aspatial or external-spatial), number of predictors, and their interaction factors. The result showed that AUC was only effected by evaluation methods ($P = 0.001$) and model algorithm × evaluation methods among these factors. In addition, the results of interaction plots (Fig. 2) showed that internal-aspatial AUC were usually larger than external-spatial AUC across all models with same variables (Fig. 3a), and for RF and TN models (Fig. 2b). However, we found only TN model had significant difference between internal-aspatial and external-spatial AUC with Paired t-test ($P = 0.000$, t=10.727, df=7), but not for RF model ($P = 0.221$, t=1.344, df=7).

<span style="color:red">Put Figure 3 Here</span>

## Prediction map assessment

According to the distribution of Hooded Crane known to us and also when compare with the source of "The BirdLife International Red Data Book" (Collar et al., 2001), for Random Forest, first, we listed RF1, RF3, RF5, RF8 as the worst predictor set of models (ranked as the fourth place; see Fig. 3 and Table 2). It shows that models with the least predictors actually perform worst. This is due to the distribution map not reflecting well the true distribution situation of Hooded Cranes, and it goes against the ecological amplitude of hooded crane distribution and biology, especially in the Far Eastern part (Fig. 4a). Second, it should be fewer areas predicted as the winter distribution area in Sakhalin Island (Russia), and whether Lake Biwa (Japan) can be predicted. Sakhalin Island is just recorded as rare breeding area of Hooded Crane in history (Collar et al., 2001). RF78, RF29 predicted slightly better than RF21 and RF11 judged on the Sakhalin Island prediction that too many areas were predicted in the middle and upper part and along this island in the RF21 and RF29 model. Third, there were largely areas of RIO ranging from 0.41 to 0.60 in the prediction map of RF78. Therefore, we regarded RF29 as the best and rank it as the first place in Random Forest.

For TreeNet: first, TN1, TN3 and TN5 are ranked as worst in our set for the same reason with Random Forest. Next, ranked 2, came TN8, TN21 and TN29, because it should be fewer areas predicted as the winter distribution area in Sakhalin Island (Russia) and Shanghai (China). Third, Lake Biwa (Japan) should be predicted, Poyang Lake and Dongting Lake (China) should predict more area and, meanwhile fewer areas should be included in the east coast of Vietnam. Thus, TN78 is ranked higher than TN11. We think that the model evaluation through a prediction map assessment may still carry bias in some extent (e.g. in our study, prediction maps from models with predictor number from 11 to 78 were very close), but it is less than going purely by internal metrics.

Put Figure 4 Here

Put Table 2 Here

## Alternative data assessment

Alternative presence data from GBIF were also used to evaluate model accuracy. From the Fig. 5, we found that most Random Forest model performed good, especealy of RF11, RF21,

283   RF29 and RF78. For TreeNet, TN3 performed significantly good, and TN21, TN29, TN78 looks
284   similar. Comparing with distribution maps (Fig. 3), these results were acceptable, becauce the
285   general distribution were close and it looks like good prediction models. Therefore, it was
286   difficult to distingusih which one was better than another one. Comparing the Ralative index of
287   Occurrence (RIO), we found more record points had higher RIOs in Random Forest than in
288   TreeNet.

289      Regression analysis was used to compare the AUC value with median and mean RIO of
290   each model in Fig. 5. We found there was a significant linear regression relationship among
291   internal-aspatial and external-spatial AUC with two kinds of RIOs in Random Forest model ($P <$
292   0.03), but the $R^2$ of spatial method (mean = 0.942) were larger than internal-aspatial mehthod
293   (mean = 0.706); while in TreeNet model, the linear relationship were not so obvious ($P > 0.189$).

294                           <span style="color:red">Put Figure 5 Here</span>

295   Figure 5 (a) violin plot of Random Forest with different predictors model; (b) violin plot of
296   TreeNet with different predictors model. The thick black bar in the centre represents the
297   interquartile range, the thin black line extended from it represents the 95% confidence intervals,
298   and the white dot is the median.

## DISCUSSION

300      In this study, we have proposed two methods to obtain three evaluation criteria (AUC, the
301   Kappa statics (Kappa) and the true skill statics (TSS)), we refer to them as internal-aspatial and
302   external-spatial approaches (Fig. 2). Further, we used a prediction map based on experience
303   knowledge (Fig. 4 and Table 2). Overall, regardless of the evaluation criteria (AUC and TSS) the
304   internal-aspatial or external-spatialmetric, the AUC and TSS in these models with different
305   predictors were close to each other. In comparison, Kappa performed slightly more distinct,
306   especially in the Random Forest model (Fig. 2). In addition, we found there were obviously
307   linear relationships among three evaluation criteria, no matter of what approach was used (Table
308   1).

309      When put to a test, the results of the three-way ANOVA showed that model accuracy based
310   on AUC was only influenced by the evaluation approach (internal-aspatial or external-spatial)
311   and the interaction of the evaluated metric and model algorithms (Random Forest or TreeNet).
312   Though the internal values were larger than the external-spatial value in same models in most

313　cases (Fig. 2a and Fig. 3), we found only TN model had a significant difference between the

314　internal and the spatial AUC, but not for the RF model when using a Paired t-test. It means that

315　the effect of the internal-aspatial and external-spatial metric to evaluate model accuracy was

316　close in Random Forest; but would somewhat mislead in TreeNet. Based on a rough classifying

317　system, AUC can be interpreted as follows: $\geq 0.9$ are excellent, $0.80 \sim 0.90$ "good", $0.70 \sim 0.80$

318　"fair", $0.60 \sim 0.70$ "poor"and $0.50 \sim 0.60$ "fail" (Allouche et al., 2006). Therefore, one same

319　TreeNet model would be listed as different accurate classes models when referring to internal-

320　aspatial and external-spatial metric, which was also seen in Kappa and TSS (Fig. 2b and 2c).

321　　From both of the internal-aspatial and external-spatial metrics of AUC and TSS values, we

322　found it was difficult to tell which model was better, when the number of predictors ranged from

323　3 to 78. But spatial Kappa for Random Forest showed distinctly different among models,

324　internal-aspatial Kappa showed an inconsistent result with external-spatial Kappa and the other

325　two statistic criteria. Combining the rank of prediction maps assessed through our field

326　knowledge (Fig. 4 and Table 2), we found that the Random Forest result was consistent with the

327　external-spatial Kappa metric. We thought taking external-spatial Kappa as the criteria was the

328　best choice among above for Random Forest model in our case, but maybe do not perform well

329　for all models and species. This needs more applications and study. Also, several studies have

330　criticized the kappa statistic for being inherently dependent on prevalence and they claimed that

331　this dependency introduces bias and statistical artefacts to estimates of accuracy (Byrt et al.,

332　1993; Lantz and Nebenzahl, 1996; Manel et al., 2001; McPherson et al., 2004).

333　　The high AUC values ($> 0.85$) and the slightly difference among all 16 Hooded Crane

334　models (Fig. 2a, Supplement S4) show that all models were accurate and performed similar, as

335　values above 0.75 generally indicate an adequate model performance for most applications

336　(Pearce and Ferrier, 2000). However, we would list RF1, RF3, TN1, TN3 and TN5 as 'bad

337　models' because of the poor spatial prediction map (Fig. 4). It means that both of the internal-

338　aspatial and external-spatial AUC did not perform so well to distinguish model predictions, but

339　the prediction maps did. Therefore, we argue that model accuracy evaluation should not only be

340　based on a static number, but also should care more about models' spatial prediction as assessed

341　with external-spatial data!

342　　In this research, we evaluated how the prediction map compares in the light of the experts'

343　knowledge on species and its real winter distribution, to determine which model is more accurate

344  and reliable. We clearly agree with the thought of Fielding (2002): accuracy assessments for

345  models should have a close connection to the intended use of the model. One major gain of

346  species distribution models is to model and quantify complex ecology and to support an efficient

347  conservation management, such as conservation planning (Drew et al., 2011). It could be used

348  for instance to design reserve networks that maintain biodiversity (Guisan and Thuiller, 2005;

349  Han et al., 2018). Thus, in order to obtain a more accurate and reliable species distribution model

350  in less time and with less money, an external-spatial approach will be more meaningful than just

351  to get a numerically perfect statistical model with unproven output and assumptions. In times of

352  limited and competing, science budgets, such things really matter, and when large scales are to

353  be handled well; valid inference remains 'key'. On such scales even small decimal improvements

354  can be of major value. Wilson et al. (2005) already concluded that efforts should be directed

355  towards producing the most reliable predictions for use in conservation planning, and for

356  instance to find the reserve network that is most robust to the uncertainty in the predictions.

357      In addition, alternative occurrence data from other sources were also used to assess model

358  accuracy. The results were similar with prediction maps and statistic metrics for Random Forest

359  models, but not really for TreeNet. This also was proved through the regression analysis result

360  between internal-aspatial, and external-spatial AUC and RIOs. The alternative data used, in this

361  research, was only 'presence' data.    In the future true absence data (=species are not occurring)

362  should also be collected, though collecting such absence data remains difficult. However, in all

363  the ways we used here and also studies elsewhere, most studies people used point data (like

364  presence, pseudo-absence points) to assess distribution area accuracy (Baltensperger et al., 2013;

365  Mi et al., 2017). Whether points can stand for the area (polygon) and how much representation

366  they own should be a discussion in future study. It's a question of detection distances as assessed

367  through Distance Sampling for instance!

368      In this study, we found the accuracy of highly non-parsimonious models RF78 and TN78

369  performed very well, and those were close to models with a set of predictors reaching from 11 to

370  29 across most of evaluation criteria. That means high-dimension variables models can also

371  predict species distribution very well. In contrast we found that, so far, high-dimension variables

372  models are widely avoided by ecological researchers (few study use more than 30 variables, we

373  referred 30 papers published, see Supplement S5). The published advice was that high

374  dimensionality is unwanted, dangerous (Meisel, 1972), poorly fitted or overfitted (Harrell et al.,

375    1996). But Breiman (2001b) stated for long time differently, and recent work has shown that

376    dimensionality can be a blessing (Drew et al. 2011) and any proxy predictor can often improve

377    the prediction, specifically when dealing with large scales and when decimals mean a lot beyond

378    just template thresholds. This presents a massive paradigm shift for the sciences. Though using

379    complex predictors may be unpleasant perhaps, and requires skill and some time, the soundest

380    path for valid inference – the goal of science (generalization; as per textbook) - is to go for

381    predictive accuracy first, then try to understand why and to infer (Breiman, 2001b; Hilborn and

382    Mangel 1997, Drew et al. 2011).

## Conclusion

384    In this study, we used two methods to assess model accuracy: internal-aspatial and external-

385    spatial. We found internal-aspatial and external-spatial metrics can get higher model

386    performance (AUC > 0.85), but both of them can't distinguish models with different predictors

387    well, while the prediction maps did a little better than them. Therefore, we argued that model

388    accuracy evaluation also should care more about models' spatial prediction and has a close

389    connection to the intended use of the model! Certainly, all above conclusion is limited to

390    Random Forest and TreeNet from lots of SDM options available, and only one species. Whether

391    other algorithm implementations and species have the same results should be tested further. As

392    we know, other than Breiman (2001b) and related papers by the authors, this maybe the first time

393    the concept of external-spatial assessment criteria for model accuracy is formerly promoted, and

394    with a quest to assess model accuracy through prediction maps for inference and applications.

395    Overall, we hope to see more study on proposing better methods and data to assess species

396    distribution models (SDMs) and prediction distribution map for valid inference, and sustainable

397    conservation management worldwide.

## ACKNOWLEDGEMENTS

404  models, and Salford Systems Ltd. This research was funded by The National Forestry Bureau of
405  China.
406

# REFERENCES

Allouche, O., Tsoar, A., Kadmon, R., 2006. Assessing the accuracy of species distribution models: prevalence, kappa and the true skill statistic (TSS). *Journal of Applied Ecology* 43, 1223-1232.

Anderson, R.P., Gonzalez Jr, I., 2011. Species-specific tuning increases robustness to sampling bias in models of species distributions: an implementation with Maxent. *Ecological Modelling* 222, 2796-2811.

Aspinall, R., Veitch, N., 1993. Habitat mapping from satellite imagery and wildlife survey data using a Bayesian modeling procedure in a GIS. *Photogrammetric Engineering and Remote Sensing* 59, 537-543.

Austin, M., Nicholls, A., Margules, C.R., 1990. Measurement of the realized qualitative niche: environmental niches of five Eucalyptus species. *Ecological Monographs* 60, 161-177.

Baltensperger, A., Mullet, T., Schmid, M., Humphries, G., Kövér, L., Huettmann, F., 2013. Seasonal observations and machine-learning-based spatial model predictions for the common raven (Corvus corax) in the urban, sub-arctic environment of Fairbanks, Alaska. *Polar Biology* 36, 1587-1599.

Barry, S., Elith, J., 2006. Error and uncertainty in habitat models. *Journal of Applied Ecology* 43, 413-423.

Betts, M.G., Ganio, L., Huso, M., Som, N., Huettmann, F., Bowman, J., Wintle, B.A., 2009. Comment on "Methods to account for spatial autocorrelation in the analysis of species distributional data: a review" *Ecography* 32: 374-378.

Beyer, H., 2013. Hawth's Analysis Tools for ArcGIS version 3.27 (software).

Birdlife international, 2014. IUCN Red list for birds.

Booms, T.L., Huettmann, F., Schempf, P.F., 2010. Gyrfalcon nest distribution in Alaska based on a predictive GIS model. *Polar Biology* 33, 347-358.

Breiman, L., 2001a. Random forests. *Machine Learning* 45, 5-32.

Breiman, L., 2001b. Statistical modeling: The two cultures (with comments and a rejoinder by the author). *Statistical Science* 16, 199-231.

Byrt, T., Bishop, J., Carlin, J.B., 1993. Bias, prevalence and kappa. *Journal of Clinical Epidemiology* 46, 423-429.

437 Collar, N.J., Crosby, R., Crosby, M., 2001. Threatened birds of Asia: the BirdLife International
438     red data book. Volume 1. Cambridge: BirdLife International.

439 Craig, E., Huettmann F., 2008. Using "blackbox" algorithms such as TreeNet  and        Random
440     Forests for data-mining and for finding meaningful patterns, relationships and outliers in
441     complex ecological data: an overview, an example using golden eagle satellite data and an
442     outlook for a promising future. Chapter IV in Intelligent Data Analysis: Developing New
443     Methodologies through Pattern Discovery and Recovery (Hsiao-fan Wang, Ed.). IGI Global,
444     Hershey, PA, USA. pp 65 -83.

445 Cushman, S.A., Huettmann, F., 2010. Spatial Complexity, Informatics, and Wildlife
446     Conservation. Springer, Springer Tokyo Berlin Heidelberg New York.

447 Derrig, R., Francis, L., 2006. Distinguishing the forest from the TREES: A comparison of tree
448     based data mining methods, Casualty Actuarial Society Forum. Citeseer, pp. 1-49.

449 Drew, C.A., Wiersma, Y., Huettmann, F., 2011. Predictive species and habitat modeling in
450     landscape ecology. Springer.

451 Elith, J., Leathwick, J.R., 2009. Species distribution models: ecological explanation and
452     prediction across space and time. *Annual Review of Ecology, Evolution, and Systematics* 40,
453     677.

454 Elith, J., Phillips, S.J., Hastie, T., Dudík, M., Chee, Y.E., Yates, C.J., 2011. A statistical
455     explanation of MaxEnt for ecologists. *Diversity and Distributions* 17, 43-57.

456 Fielding, A.H., 2002. What are the appropriate characteristics of an accuracy measure.
457     *Predicting species occurrences: issues of accuracy and scale*, 271-280.

458 Fielding, A.H., Bell, J.F., 1997. A review of methods for the assessment of prediction errors in
459     conservation presence/absence models. *Environmental conservation* 24, 38-49.

460 Franklin, J., 1995. Predictive vegetation mapping: geographic modelling of biospatial patterns in
461     relation to environmental gradients. *Progress in Physical Geography* 19, 474-499.

462 Franklin, J., McCullough, P., Gray, C., 2000. Terrain variables used for predictive mapping of
463     vegetation communities in Southern California. Terrain analysis: principles and
464     applications/edited by John P. Wilson, John C. Gallant.

465 Fernandez-Delgado, M., Cernadas, E, Barro, S., Dinan, A., 2014. Do we Need Hundreds of
466     Classifiers to Solve Real World Classification Problems? *Journal of Machine Learning* 15:
467     3133-3181.

468    Friedman, J., 2001. TreeNet™.

469    Friedman, J.H., 2002. Stochastic gradient boosting. *Computational Statistics and Data Analysis*
470        38, 367-378.

471    Guisan, A., Lehmann, A., Ferrier, S., Austin, M., OVERTON, J., Aspinall, R., Hastie, T., 2006.
472        Making better biogeographical predictions of species' distributions. *Journal of Applied*
473        *Ecology* 43, 386-392.

474    Guisan, A., Thuiller, W., 2005. Predicting species distribution: offering more than simple habitat
475        models. *Ecology letters* 8, 993-1009.

476    Guisan, A., Zimmermann, N.E., 2000. Predictive habitat distribution models in ecology.
477        *Ecological modelling* 135, 147-186.

478    Guo, Y., 2005. Breeding behaviror ecology of Hooded Crane (Grus monacha) in Lesser Xingan
479        Mountains. Northeast Forest University.

480    Han, X., Huettmann, F., Guo, Y., Mi, C., Wen, L., 2018. Conservation prioritization with
481        machine learning predictions for the black-necked crane Grus nigricollis, a flagship species
482        on the Tibetan Plateau for 2070. *Regional Environmental Change* 1-10.

483    Harrell, F., Lee, K.L., Mark, D.B., 1996. Tutorial in biostatistics multivariable prognostic models:
484        issues in developing models, evaluating assumptions and adequacy, and measuring and
485        reducing errors. *Statistics in medicine* 15, 361-387.

486    Hastie, T., Tibshirani, R., Friedman, J., 2009. The Elements of Statistical Learning: Data Mining,
487        Inference, and Prediction. Second Edition.Springer New York.

488    Hegel, T.M., Cushman, S.A., Evans, J., Huettmann, F., 2010. Current state of the art for
489        statistical modelling of species distributions, Spatial complexity, informatics, and wildlife
490        conservation. Springer, pp. 273-311.

491    Herrick, K.A., Huettmann, F., Lindgren, M.A., 2013. A global model of avian influenza
492        prediction in wild birds: the importance of northern regions. *Vet Res* 44, 42.

493    Hilborn, R., Mangel, M.,1997. The Ecological Detective: Confronting Models with Data. 1997.
494        Princeton New Jersey.

495    Huettmann, F., Gottschalk, T., 2011. Simplicity, Model Fit, Complexity and Uncertainty in
496        Spatial Prediction Models Applied Over Time: We Are Quite Sure, Aren't We?, Predictive
497        Species and Habitat Modeling in Landscape Ecology. Springer, pp. 189-208.

498    IUCN, 2012. IUCN red list of threatened species. Version 2012.2. International Union for the

499       Conservation of Nature Gland, Switzerland.

500    Jiao, S., Huettmann, F., Guo, Y., Li, X., Ouyang, Y., 2016. Advanced long-term bird banding

501       and climate data mining in spring confirm passerine population declines for   the  Northeast

502       Chinese-Russian flyway. *Global and Planetary Change* 144: 17-33

503    Kadmon, R., Farber, O., Danin, A., 2003. A systematic analysis of factors affecting the

504       performance of climatic envelope models. *Ecological Applications* 13, 853-867.

505    Kandel K. Huettmann, F., Suwal, M. K., Regmi, G.R., Nijman, V ., Nekaris, K.A.I.,

506       Lama, S.T., Thapa, A., Sharma, H.P., Subedi, T.R., 2015. Rapid multi-nation distribution

507       assessment of a charismatic conservation species using open access ensemble model

508       GIS predictions: Red panda (Ailurus fulgens) in the Hindu-Kush Himalaya region.

509       *Biological Conservation* 181: 150-161.

510    Lantz, C.A., Nebenzahl, E., 1996. Behavior and interpretation of the statistic: Resolution of the

511       two paradoxes. *Journal of clinical epidemiology* 49, 431-434.

512    Loiselle, B.A., Howell, C.A., Graham, C.H., Goerck, J.M., Brooks, T., Smith, K.G., Williams,

513       P.H., 2003. Avoiding pitfalls of using species distribution models in conservation planning.

514       *Conservation biology* 17, 1591-1600.

515    Manel, S., Williams, H.C., Ormerod, S.J., 2001. Evaluating presence–absence models in ecology:

516       the need to account for prevalence. *Journal of applied Ecology* 38, 921-931.

517    McPherson, J., Jetz, W., Rogers, D.J., 2004. The effects of species' range sizes on the accuracy

518       of distribution models: ecological phenomenon or statistical artefact? *Journal of applied*

519       *ecology* 41, 811-823.

520    Mi, C., Huettmann, F., Guo, Y., 2014. Obtaining the best possible predictions of habitat selection

521       for wintering Great Bustards in Cangzhou, Hebei Province with rapid machine learning

522       ßanalysis. *Chinese Science Bulletin*, 1-9.

523    Mi, C., Huettmann, F., Guo, Y., 2016. Climate envelope predictions indicate an enlarged suitable

524       wintering distribution for Great Bustards (Otis tarda dybowskii) in China for the 21st

525       century. *PeerJ* 4:e1630.

526    Mi, C., Huettman, F., Guo, Y., Han, X., Wen, L., 2017. Why choose Random Forest to predict

527       rare species distribution with few samples in large undersampled areas? Three Asian crane

528       species models provide supporting evidence. *PeerJ* 5:e2849.

529    Mi, C., Møller, A.P., Guo, Y., (2018) Annual spatio-temporal migration patterns of Hooded

Cranes wintering in Izumi based on satellite tracking and their implications for conservation. *Avian Research* 9, 23.

Oppel, S., Meirinho, A., Ramírez, I., Gardner, B., O'Connell, A.F., Miller, P.I., Louzao, M., 2012. Comparison of five modelling techniques to predict the spatial distribution and abundance of seabirds. *Biological Conservation* 156, 94-104.

Pearce, J., Ferrier, S., 2000. Evaluating the predictive performance of habitat models developed using logistic regression. *Ecological modelling* 133, 225-245.

Pearson, R.G., Dawson, T.P., Liu, C., 2004. Modelling species distributions in Britain: a hierarchical integration of climate and land-cover data. *Ecography* 27, 285-298.

Pearson, R.G., Thuiller, W., Araújo, M.B., Martinez-Meyer, E., Brotons, L., McClean, C., Miles, L., Segurado, P., Dawson, T.P., Lees, D.C., 2006. Model-based uncertainty in species range prediction. *Journal of Biogeography* 33, 1704-1711.

Peterson, A.T., Ortega-Huerta, M.A., Bartley, J., Sánchez-Cordero, V., Soberón, J., Buddemeier, R.H., Stockwell, D.R., 2002. Future projections for Mexican faunas under global climate change scenarios. *Nature* 416, 626-629.

Prasad, A.M., Iverson, L.R., Liaw, A., 2006. Newer classification and regression tree techniques: bagging and random forests for ecological prediction. *Ecosystems* 9, 181-199.

Reese, G.C., Wilson, K.R., Hoeting, J.A., Flather, C.H., 2005. Factors affecting species distribution predictions: a simulation modeling experiment. *Ecological Applications* 15, 554-564.

Regmi G.R., Huettmann, F., Suwal, M. K., Nijman, V., Nekaris, K.A.I., Kandel, K., Sharma, N. and Coudrat, C., 2018. First Open Access Ensemble Climate Envelope Predictions of Assamese Macaque Macaca Assamensis in South and South-East Asia: A new role model and assessment of endangered species. *Endangered Species Research* 36: 149–160.

Segurado, P., Araujo, M.B., 2004. An evaluation of methods for modelling species distributions. *Journal of Biogeography* 31, 1555-1568.

Seoane, J., Carrascal, L.M., Alonso, C.L., Palomino, D., 2005. Species-specific traits associated to prediction errors in bird habitat suitability modelling. *Ecological Modelling* 185, 299-308.

Simonov, E.A., Dahmer, T.D., 2008. Amur-Heilong river basin reader. Ecosystems.

Vaughan, I., Ormerod, S., 2005. The continuing challenges of testing species distribution models. *Journal of Applied Ecology* 42, 720-730.

561     Wei, C.L., Rowe, G.T., Escobar-Briones, E., Boetius, A., Soltwedel, T., Caley, M.J., Soliman, Y.,

562        Huettmann, F., Qu, F., Yu, Z., 2010. Global patterns and predictions of seafloor biomass

563        using random forests. *PLoS One* 5, e15323.

564     Wilson, K.A., Westphal, M.I., Possingham, H.P., Elith, J., 2005. Sensitivity of conservation

565        planning to different approaches to using predicted species distribution data. *Biological*

566        *Conservation* 122, 99-112.

567

568  Figure Legend

569  **Figure 1** Map of study area and study species locations.

570

571  **Figure 2** (a) barplot of AUC from internal-aspatial and external-spatial metrics of RF and TN; (b)

572  barplot of Kappa from internal-aspatial and external-spatial source of RF and TN; (c) barplot of

573  TSS from internal-aspatial and external-spatial source of RF and TN.

574

575  **Figure 3** Interaction plot of AUC value between predictor number, model algorithms and

576  evaluation methods. (a) Interaction plot of AUC value between predictor number (1, 3, 5, 8, 11,

577  21, 29 and 78) and two evaluation methods (internal-aspatial and external-spatial), (b) Interaction

578  plot of AUC between model algorithms (Random Forest) and two evaluation methods (internal-

579  aspatial and external-spatial). Dots with same line represent the AUC value (internal-aspatial and

580  external-spatial) of same model with certain predictor number (1, 3, 5, 8, 11, 21, 29 and 78).

581

582  **Figure 4** Prediction map of Random Forest and TreeNet with 8 different predictor numbers. (a)

583  Prediction map of Random Forest; (b) Prediction map of TreeNet.

584

585  **Figure 5** (a) violin plot of Random Forest with different predictors model; (b) violin plot of

586  TreeNet with different predictors model. The thick black bar in the centre represents the

587  interquartile range, the thin black line extended from it represents the 95% confidence intervals,

588  and the white dot is the median.

589

590

591

592

593

594

595

596

597

598

599

600     Table Legend

601

602     **Table 1** Linear Regression analysis of AUC, Kappa and TSS.

603

604     **Table 2** Rank of model by prediction map assessment.

605

# Table 1(on next page)

Table

Table

1          Table 1 Linear Regression analysis of AUC, Kappa and TSS

|  | Slope | R² | *P* |
|---|---|---|---|
| AUC~Kappa | 0.165 | 0.510 | 0.000 |
| AUC~TSS | 0.481 | 0.839 | 0.000 |
| AUC_RF~Kappa_RF | 0.159 | 0.556 | 0.001 |
| AUC_RF~TSS_RF | 0.455 | 0.773 | 0.000 |
| AUC_TN~Kappa_TN | 0.203 | 0.582 | 0.000 |
| AUC_TN~TSS_TN | 0.533 | 0.937 | 0.000 |
| AUC_RF_spatial~Kappa_RF_spa | 0.203 | 0.847 | 0.001 |
| AUC_RF_spa~TSS_RF_spa | 0.507 | 0.999 | 0.000 |
| AUC_RF_aspa~Kappa_RF_aspa | 0.120 | 0.348 | 0.123 |
| AUC_RF_aspa~TSS_RF_aspa | 0.598 | 0.970 | 0.000 |
| AUC_TN_spa~Kappa_TN_spa | 0.234 | 0.968 | 0.000 |
| AUC_TN_spa~TSS_TN_spa | 0.503 | 1.000 | 0.000 |
| AUC_TN_aspa~Kappa_TN_aspa | 0.074 | 0.968 | 0.000 |
| AUC_TN_aspa~TSS_TN_aspa | 0.263 | 0.943 | 0.000 |

2    Note: aspa represents internal-aspatial metric, spa represents external-spatial metric.

3

4

Table 2 Rank of model by prediction map assessment.

| Rank | Random Forest | TreeNet |
|------|---------------|---------|
| 1 | RF29 | TN78 |
| 2 | RF78 | TN11, |
| 3 | RF11, RF21 | TN29, TN21, TN8 |
| 4 | RF1, RF3, RF5, RF8 | TN1, TN3, TN5 |

5    Note: 1 means the best; 2 means better; 3 means good; 4 means less good.

**Figure 1**(on next page)

Figure 1 Study area

**Figure 2**(on next page)

Figure 2

**Figure 3**(on next page)

Figure 3

**Figure 4**(on next page)

Figure 4a

(a)

1 predictor

3 predictors

5 predictors

8 predictors

11 predictors

21 predictors

29 predictors

78 predictors

**Legend**
Random Forest
with different predictors

Predicted Occurrence

| | |
|---|---|
| | 0.00 - 0.20 |
| | 0.21 - 0.40 |
| | 0.41 - 0.60 |
| | 0.61 - 0.80 |
| | 0.81 - 1.00 |

**Figure 5**(on next page)

Figure 4b

(b)

1 predictor | 3 predictors

5 predictors | 8 predictors

11 predictors | 21 predictors

29 predictors | 78 predictors

**Legend**

DooNet with
different predictors

**Predicted Occurrence**

- 0.00 - 0.20
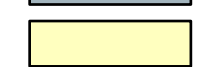- 0.21 - 0.40
- 0.41 - 0.60
- 0.61 - 0.80
- 0.81 - 1.00

**Figure 6**(on next page)

Figure 5