**Archimer**
https://archimer.ifremer.fr

# Which spatial interpolators I should use? A case study applying to marine species

Rufino Marta [1, 2, 3, *], Albouy Camille [1], Brind'Amour Anik [1]

[1] IFREMER - Centre Atlantique, French Research Institute for Exploitation of the Sea, Département Ecologie et Modèles pour l'Halieutique (EMH), Rue de l'Ile d'Yeu - BP 21105, 44311 Nantes cedex 3, France
[2] Portuguese Institute for the Sea and the Atmosphere (IPMA), Division of Modelling and Management of Fisheries Resources, Av. Dr. Alfredo Magalhães Ramalho, 6, 1495-165 Lisboa, Portugal
[3] Centre of Statistics and its Applications (CEAUL), Faculty of Sciences, University of Lisbon, Portugal

* Corresponding author : Marta Rufino, email address :  marta.rufino@ipma.pt

**Abstract :**

Species are spread in space, whereas sampling is sparse. Thus, to describe and map along environmental gradients, it is necessary to interpolate the species abundance. Considering the plethora of valid methods, the researcher gets easily puzzled to choose the most appropriate interpolation approach with reference to the ecological question being asked.
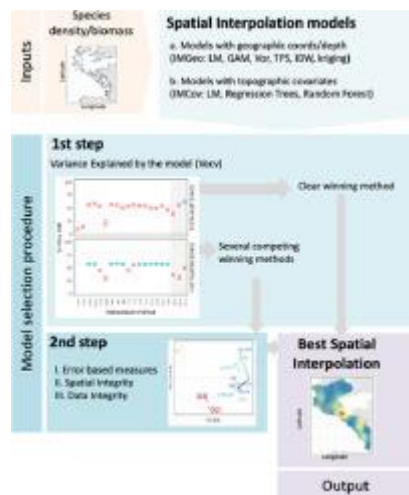
We propose a procedure to select among alternative spatial distribution models and we illustrate it with 175 marine species distributions (35 species * 5 years). In a first step, the distribution of the variance explained by the predictive model (VEcv) given by 10-fold cross validation is estimated for each interpolation method. When the inter-quartile range of the VEcv distribution of the different methods overlap, the selection passes to a second step, using 11 measures belonging to three criteria: 1) error based measures, 2) spatial equivalence measures (center of gravity, inertia, isotropy and index of aggregation) and 3) measures based on the data integrity after interpolation, for example the percentage of area over the maximum sampled data.

We applied our approach to marine species sampled using either stratified random survey (trawl) or systematic survey (acoustic). We found that 87% of all species distributions had overlapping VEcv and thus passed the first selection. In the second selection step, the best method varied with species and year, although general additive model (GAM), Thin Plate Spline (TPS), Universal Kriging (UKr) and Random Forest (Rfor) performed better for the trawl data and TPS, Ordinary Kriging (OKri) and UKr for the acoustic data. Further, the results differed within methods (e.g. kriging neighborhood and type of kriging) and small modifications on the specifications can have a large impact on the surfaces produced.

The proposed approach 1) is accessible and intuitive, and does not require any complex software or sophisticated methodology; 2) shows exactly in what aspects each interpolation model is prevalent over the others and permits to make a decision accordingly to the objectives of the study; 3) takes into account

different criteria to evaluate each, properties of an interpolation method; 4) is universal and does not depend on the method used or the data characteristics. A detailed review on the subject is also included.

**Graphical abstract**



**Highlights**

► A new method was developed to select among spatial distribution models using 2-steps. ► The 1st step uses the variance explained by the predictive model (10-fold cross validation). ► The 2nd step uses 3 criteria: error based measures, spatial equivalence measures and data integrity. ► The method is illustrated using 175 marine species distributions (35 species x 5. years). ► The approach is accessible, clear, multi-criteria and is universal as it does not depend on the method used or the data characteristics.

## Introduction

Everything becomes clearer when it is presented in a figure. Maps of species distributions are required for numerous purposes, such as to visualize spatial variability, spot changes in the communities or provide estimations of variables of interest. Species distributions are mapped by interpolation, that is, predicting values at un-sampled areas using a modelling procedure applied to the sampled data. Ecological processes are inherent in species distributions (or spatial structure), such as those impacted by anthropogenic factors and climate change, and thus reflected in the respective maps. Coupled with species traits or phylogenetic information, species distribution maps can inform on the location of functional or phylogenetic hotspots.

Interpolated species distributions are widely used in a range of fields and applications, including regional biodiversity assessments, spatial conservation prioritization, evolutionary biology, epidemiology, global change biology and wildlife management (Araujo and Peterson, 2012). Given the importance of spatial interpolation, new and more powerful methods are developed on a regular basis and continuous/progressive evaluation of these statistical models is necessary (Austin, 2007).

There is a large body of literature on species distribution models (SDM, also known as bioclimatic envelope models, ecological niche models and habitat suitability models), that explore the relationship between geographical occurrences of species and corresponding environmental variables (Araújo and Guisan, 2006; Dormann et al., 2007; Elith and Leathwick, 2009; Guisan and Zimmermann, 2000; Hui et al., 2013; Olden and Jackson, 2002). These are however, more challenging to apply on marine data to model than on its terrestrial counterparts. Marine fisheries data compared to its terrestrial counterpart (Lecours et al., 2016), typically have fewer sampling stations, in face of the large costs associated with the survey operations and cover irregular survey shape areas (e.g. 3-15 m depth along the coast). Further, marine species distributions are typically characterized by a large percentage of zero observations and a huge variability. Finally, the availability of environmental covariates in the marine field is generally scarcer, both in terms of geographic span and resolution, also due to difficulties associated with sampling. Thus, in face of this particularities of marine data, the challenges arisen by producing spatial models using this type of data have been often referred (Olden and Jackson, 2002).

So, the next question that naturally arises is which method to choose? It is widely recognized that there are no magical recipes to determine the perfect interpolation model. Undoubtedly, this should be focused on the data *per se*, laying on the statistical theoretical background and respective assumptions (Bivand et al., 2013; Cressie, 1993; Li and Heap, 2008; Sluiter, 2009; Wackernagel, 1998; Webster and Oliver, 2007). Several methods pass this first selection, however we need a framework accounting for model selection and evaluation to help decision making. Such decision is a complex issue, central to ecological modelling, with huge implications (Naimi and Araújo, 2016). Model selection involves evaluation, validation, performance, accuracy, skill, efficiency or robustness. Nevertheless these concepts are hard to disentangle and have been used with different meanings on the literature (reviewed in Bellocchi et al., 2010).

Generally, interpolators are compared using error based measurements, that is predicted *vs*. observed values, preferably obtained from cross-validation or jackknife processes (Li, 2016; reviewed by Li and Heap, 2008, 2011; Richter et al., 2012; Stow et al., 2009; Willmott et al., 2015). Using primarily error-based criteria, comparisons among spatial interpolators have been done in the field of meteorology (Aalto et al., 2013), air quality (Hoffman, 2015), soil (Gasch et al., 2015; Hengl et al., 2004, 2015), marine sediment (Diesing et al., 2014; Lark et al., 2016; Li et al., 2011), bathymetry (Amante and Eakins, 2016) and environmental sciences (Li and Heap, 2011, 2008). Among these, Mean Absolute Error (MAE, bias or a measure of average

103 error-magnitude) and Root Mean Squared Error (RMSE, for accuracy or over/under fitting) are
104 the most commonly within the field of environmental sciences (Li and Heap, 2008; Richter et
105 al., 2012; Willmott, 1982). However, as their magnitude depends on the scale/unit of the
106 variable predicted, these are hardly comparable among variables or subjects. Further, most of
107 these accuracy measures are algebraically related, being thus potentially redundant and
108 collinear (Li, 2016; Li and Heap, 2011, 2008; Willmott et al., 2015). It has also been suggested
109 the use of dimensionless measures, besides at least one error measure on the variable scale
110 (modified coefficient of efficiency, Legates and McCabe Jr., 1999; coefficient of efficiency,
111 Nash and Sutcliffe, 1970; index of aggreement, Willmott, 1981, 1982; modified index of
112 aggreement, Willmott et al., 2012, 2015). Such approaches have been applied within the field
113 of hydrology/climatology until recently, where Li (2016, 2017) revising Willmot's D,
114 advocated its use as an universal tool to assess the accuracy of predictive models within
115 environmental sciences, naming it Variance Explained by predictive models, estimated by
116 cross-validation (VEcv), that is: how well a model is predicted, relative to the average of the
117 observations (also called coefficient of efficiency, Nash and Sutcliffe, 1970 or G-value or
118 goodness-of- prediction measure).
119 However, all these criteria often lead to overlapping results, that is several models having
120 similar VEcv values, being difficult to select only one. Furthermore, these are based essentially
121 on error measurements and do not take into account other fundamental aspects, such as spatial
122 integrity (that is if the interpolation respects the spatial distribution of the data) or the spatial
123 data limits of the interpolation relative to the original data. Additionally, whatever the criteria
124 considered, it has long been demanded the establishment of a consistent and rationale set of
125 procedures that should be used to compare spatial interpolation models (Fox 1981) (Willmott,
126 1982).
127 The objective of the current work, is therefore, to develop a simple and accessible protocol to
128 compare the results given by different spatial interpolating methods, that integrates important
129 aspects for mapping marine species distribution, namely not only error measures, but also
130 spatial and data integrity after interpolation. The proposed protocol was applied to compare 20
131 interpolation methods applied to 35 species distributions from two typical fisheries surveys
132 (trawl and acoustic), carried out during 5 years. The interpolation methods considered comprise
133 approaches only using geographic coordinates and methods using depth and other topographic
134 variables derived from bathymetry.

## Materials and Methods

136 DATA

137 We considered two data sets for the species distributions case studies, one obtained from
138 scientific groundfish bottom trawl surveys (EVHOE) and another from a scientific pelagic
139 acoustic surveys (PELGAS).
140 The bottom trawl survey is carried out annually during Autumn in the North Atlantic
141 ("Evaluation Halieutique de l'Ouest Européen, EVHOE cruise, RV Thalassa, IFREMER,"
142 n.d.). It ranges from the Bay of Biscay up to the Celtic sea, with a randomly stratified sampling
143 strategy, comprising from 119 to 153 stations/year ("Evaluation Halieutique de l'Ouest
144 Européen, EVHOE cruise, RV Thalassa, IFREMER," n.d.)(map and location of sampling
145 stations can be found in see Supplement 1). The biomass of the 29 fish species occurring more
146 than 10 times/year during the survey and excluding the main pelagic species was used (see
147 Supplement 1 for the species list names and further details on the survey). For the purpose of

148    this study we used data between 2011 and 2015 with an average of 198 sampling stations
149    (number of hauls per year can be found in Supplement 1).

150    The pelagic survey (Doray, M., Duhamel E. , Huret M. , Petitgas P., 2002; Doray M., Badts V.,
151    Masse J., Duhamel E., Huret M., Doremus G., 2014) is an acoustic spring survey that aims at
152    monitoring the Bay of Biscay pelagic ecosystem to inform fisheries and ecosystem
153    management. Initially, PELGAS objective was to estimate biomass anchovy (*Engraulis*
154    *encrasicolus*) and nowaday, the survey goals were extended to estimate the stocks of all the
155    small pelagic fish species in the Bay of Biscay. From this survey we extracted the biomasses of
156    six small pelagic fish species (see Supplement 1 for the species list and map), sampled over
157    1345 to 1997 locations obtained from 29 acoustic radials perpendicular to the coast between
158    2011 and 2015.

159    In the simplest cases, interpolation can be carried out using only the geographic coordinates as
160    explanatory variables. However, in marine systems bathymetry influences species spatial
161    distribution and this information is available at global scale, and thus can be added to improve
162    interpolation models. Bathymetry data was extracted from GEBCO data base and validated
163    with depth data obtained during the EVHOE surveys (IMGeo). Additionally, other covariates
164    derived from bathymetry can be included to improve interpolation models. Those additional
165    covariates were added without referring to any specific ecological hypotheses, but likely
166    serving as proxy of other unmeasured environmental variables. We extrapolated eleven
167    covariates (IMCov) such as derived from bathymetry/elevation such as (Lecours et al., 2016;
168    Wilson et al., 2007): slope, aspect, northerness, easterness, rough, profile curvature (surf.curv),
169    bathymetric position index (TPI), terrain ruggedness index (TRI), surface flow (flowdir), local
170    Moran I (moran) and distance to the nearest coast (dist.coast)(details of each variable and
171    respective maps can be found in Supplement 2). Thus, all models using covariates (IMCov)
172    were produced including all variables, whereas the final model, was produced with an automatic
173    selection of these variables for each distribution.

174    INTERPOLATION METHODS

175    Seven families of methods were applied to both case studies, aiming not be exhaustive: linear
176    models (LM), general additive models (GAM), Inverse Distance Weighting (IDW), Thin Plate
177    Spline (TPS), VORonoi triangulation (VOR), Kriging (Kr) and stochastic Conditional
178    Simulation (CSim) for the methods using just geographic coordinates and eventually depth
179    (dep)(IMGeo) and three families of methods using the 11 bottom topographical variables
180    (IMCov): multiple regression (GLM), Regression Tree (RTre) and Random Forest (RFor)
181    (Hengl et al., 2015, 2007, 2004, 2003; Li et al., 2011)(Supplement 3). Additionally, we used
182    several alternatives on some methods, intended to quantify the within-model and between-
183    model variability (Araújo and Guisan, 2006), namely considering only geographic variables as
184    covariates or adding depth as well (GAMl/ GAMd, MKri/ UKri), changing kriging
185    neighborhood (OK03/ OK05/ OK07/ OK10/ OK20/ OK30 or the fitting procedure as automatic
186    vs. manual (OKri/MKri)). We provided a brief description of the methods used in Supplement
187    3 whereas additional details can be found in the vast literature (Bivand et al., 2013; Cressie,
188    1993; Fortin and Dale, 2005) and more precisely in two reviews on the subject (Li and Heap,
189    2008; Sluiter, 2009).

190 CRITERIA FOR COMPARISON OF INTERPOLATORS

191 Three complementary criteria were used to compare and evaluate the accuracy of interpolation
192 models: (i) error based measures, (ii) changes in the spatial structure due to interpolation and
193 (iii) data integrity after interpolation.

194     *1.  Error based measures*

195 Error indices were estimated using predicted and observed values obtained by ten-fold cross-
196 validation (10-fold CV). Ten-fold cross-validation was done by randomly splitting the data into
197 10 parts. We estimated the model using 9 of those 10 data set, whereas the observed values
198 ($10^{th}$ split) are predicted using the model estimated. This process is repeated for each of the ten
199 splits, obtaining predicted and observed values for the ten folds, which are then used to estimate
200 the error measures. We performed a 10-fold CV, and instead of leave one out procedure, as this
201 method has been considered to give too optimistic measures of error. The process of 10-fold
202 CV was then repeated with a random split 100 times to obtain a distribution of the error indices.
203 From predicted and observed values obtained by 10-fold, three measures were calculated: MAE
204 (Mean Absolute Error), RMSE (Root Mean Squared Error) and the Variance explained by
205 predicted models estimated using cross validation procedures (VEcv). For comparability with
206 previous works, the MAE ( $[0, \infty]$, the lower the better) and Root Mean Squared Error (RMSE,
207 $[0, \infty]$, the lower the better) were estimated (Richter et al., 2012).

208
$$MAE = \frac{1}{n}\sum_{i=1}^{n}(|P_i - O_i|)$$

209
$$RMSE = \sqrt{\frac{1}{n}\sum_{i=1}^{n}(P_i - O_i)^2}$$

210 VEcv a dimensionless measure, varies between 100 for an excellent model and $-\infty$. VEcv lower
211 than 0, indicates that the model is worse than the average and we rounded these values to -1 for
212 visualization purposes.

213
$$VEcv = 100 \times (1 - \frac{\sum(P_i - O_i)^2}{\sum(O_i - \bar{O})^2})$$

214 To perform the first step of our selection procedure (Fig. 1), we estimated the average RMSE,
215 MAE and VEcv from the distributions obtained by the 10-fold CV, along with the upper and
216 lower quantiles of the VEcv (probability of 0.25 and 0.75). We classified the average VEcv
217 following Li (2016), into: 1) very poor if VEcv ≤ 10%; 2), poor if 10 < VEcv ≤ 30%; 3), average
218 30 < VEcv ≤ 50%; 4), good if 50 < VEcv ≤ 80%; 5) and excellent if VEcv > 80%.
219 To harmonize the interpolation with the other interpolators considered in this study (MAE,
220 RMSE), we calculated an inversion of the VEcv (VEcv.inv = abs(VEcv/100-1)*100), where
221 the lowest value, indicates the worst the model.

222     *2.  Spatial integrity*

223 Ideally, an interpolation method should preserve the geometrical properties of the data, and
224 therefore its spatial structure. We evaluated changes in the data spatial structure due to
225 interpolation using deviations on four spatial indicators, relative to its results on the sampled
226 data. Four spatial indicators were considered in the current work, following the revision and

227 work of (Rufino et al., 2019). The center of gravity (CG) indicates the mean spatial location of
228 the population (Bez and Rivoirard, 2001; Woillez et al., 2009b). The Euclidean distance
229 between the CG estimated using the sampled data (without considering different areas of
230 influence for each sample) and the interpolated data was calculated. It was used as a measure
231 of the impact of interpolation on the geometric center of the distribution. The inertia represents
232 the spatial dispersion of the population around its CG, i.e. the mean square distance between
233 individual fish and the CG (Bez and Rivoirard, 2001; Woillez et al., 2009b). Isotropy/anisotropy
234 (isotropy) represents the dispersion shape of the inertia around the CG (i.e. round or ellipsoid),
235 and it is simply the ratio between the two inertia axes (Woillez et al., 2009a).
236 The Gini index quantifies the distribution's aggregation or concentration, and represents twice
237 the area between the identity function and the Lorenz curve. It is bounded between 1 and 0, and
238 the highest its value the most concentrated is the biomass in fewest samples.
239 To quantify geometric changes in the spatial distribution due to interpolation, the difference
240 between the spatial indicators calculated for the sampled data and the indicators calculated in
241 from the interpolated surfaces was estimated Thus, the absolute difference between inertia (log
242 transformed for scaling), isotropy and Gini index of the interpolated surfaces and the respective
243 ones estimated on the sampled data ($\Delta I = | I_{sample} - I_{interpolation}|$) were used as a measure of the
244 interpolation method spatial integrity. Therefore, the higher the value of the difference on these
245 spatial measures, the higher the shift relative to the sampled data's spatial structure cause by
246 the interpolation method, thus the worse its performance.

247 ### 3. Data limits integrity

248 Ideally, an interpolation method should preserve the samples data integrity and not predict
249 values outside the data range or shifts in the mean biomass. To evaluate interpolator's data
250 integrity, we used four measures estimated on each interpolated surface:
251 1) a.pix_under: relative area below/above minimum sampled biomass (pix_under =
252 $abs[(\#pixels[max(B)_{interpolation} < max(B)_{sample}]/ \#pixels_{interpolation}]\times100)$;
253 2) a.pix_over: relative area above/bellow maximum sampled prediction (pix_over)
254 $(abspix\_over = [(\#pixels[max(B)_{interpolation} > max(B)_{sample}]/ \#pixels_{interpolation}]\times100)$;
255 3) a.mean_perc: relative change on the predicted mean biomass (mean_perc = $abs[(\mu_{interpolation}$
256 $- \mu_{sample})/\mu_{sample}]\times100)$;
257 4) a.over_perc: relative change on the maximum prediction (over_prec = $abs[(\mu_{interpolation} -$
258 $\mu_{sample})/\mu_{sample}]\times100)$;
259 Where $B_{interpolation}$ is the interpolated biomass, $B_{sample}$ is the samped biomass and #pixels, the
260 number of pixels. As most of these measures were dimensionless, they can be compared among
261 studies and methods and interpret as follow: the higher the value, the wose the performance of
262 the interpolator (for comparability with the other indicators).
263 The eleven criteria were scaled such as the highest the value, the greatest the impact of the
264 interpolation on the spatial structure and data integrity or greater the errors, bias or accuracy.
265 Additionally, these criteria were not strongly correlated between each other, except for the
266 RMSE/MAE (which is expected algebraically), a.mean_perc with a.over_perc (r=-63%) and
267 CGdist with a.inertia (r = 58%). Further, these criteria were also not strongly correlated with
268 biomass, which is another desirable property. The correlation between indicators is shown in
269 supplement 4.

PROTOCOL TO COMPARE THE INTERPOLATORS

271 The 2-steps protocol proposed to compare spatial interpolators is summarized in Fig. 1.
272 First step:
273 a) <u>Estimate</u> the distribution of VEcv obtained by 10-fold CV. Select the interpolation method
274 with the highest average VEcv;
275 b) Calculate the interquartile range (IQR) of each VEcv distribution (Q1, 25% and Q3, 75%);
276 c) Select all interpolation methods where the IQR overlaps with the interpolation method that
277 showed highest VEcv;
278 d) If no other interpolation method IQR overlaps the highest VEcv IQR, the decision is reached,
279 and the best method is clear. Otherwise, we continue to the second step.
280  Second step:

281    a) Calculate all measures proposed above for each interpolation method: 3 error

282       measures (VEcv, MAE and RMSE), 4 spatial indicators (distance of the center of

283       gravity, difference in inertia, isotropy and Gini index) and the 4 data integrity

284       measures;

285    b) Do a principal component analysis (PCA), scaled and not centered on the indicators

286       matrix for the distribution (species-year), that render the criteria ranges comparable,

287       and that integrates all measures vs. all methods being evaluated;

288    c) The inverse of the distance between the methods loadings on the first two PCA axes

289       relative to the center of the PCA, is then used to rank the interpolation methods; The

290       PCA further shows concretely what aspects of the sampled data, the interpolation

291       model is not respecting;

292 This protocol was applied to every species/year distribution, for all interpolation methods being
293 assessed as a case study. All analyses were carried out using r-project. An R-script with a small
294 simplified example is added in the supplement 5.

295 **Results**

296 FIRST SELECTION STEP: VECV CRITERIA

297 Twenty two percent of the models performed very poorly, producing interpolated surfaces that
298 were worse than the mean, as shown by the negative average of the variance explained by the
299 predictive model (VEcv < 0 in 784 models; Fig. 2). The maximum variance explained by the
300 predictive model of the interpolations (VEcv) was 86% (EVHOE.TRISESM.2013-Kr3 -

301 *Trisopterus esmarkii* in 2013). Three species/years showed negative VEcv for all models, and
302 thus were eliminated from further analysis (CONGCON.2011, 2012 - *Conger conger* and
303 LOPHPIS.2014 - *Lophius piscatorius* from 2014).
304 The percentage of 'bad' models was higher for the bottom trawl survey data (EVHOE; 25%)
305 than for the pelagic survey (PELGAS; 11% Fig. 2 and Fig. 3). For the bottom trawl survey,
306 VEcv was slightly higher for GAM, TPS, geostatistical models (Kri) and RFor. As for the
307 pelagic surveys, IDW, TPS, geostatistical models (Kri) and RFor reach better results (Fig. 2
308 and Fig. 3). However, overall good VEcv were lower than 50%, except for the acoustic surveys.
309 Only few species distributions sampled with the acoustic surveys attained 'excellent' VEcv
310 classes, whereas most distributions were classified as good for IDW, TPS, Kriging and RForest
311 (Fig. 2). For the bottom trawl data, most models were classified as average or poor according
312 to the VEcv criteria (Fig. 2). LM1 and LM2 were systematically worse than the mean.
313 Overall, the use of bottom covariates did not improve the interpolation's VEcv (Fig. 3). The
314 variance explained by the predictive model (VEcv) varied more among species than between
315 years (Fig. 4 and Fig. 5). However, for some benthic species the models using covariates
316 showed higher VEcv than models without covariates (e.g. CONGCON - *Conger conger*,
317 PHYBLE- *Phycis blennoides*, SOLESOL - *Solea solea*, HELIDAC - *Helicolenus*
318 *dactylopterus*, SCYOCAN - *Scyliorhinus canicula* and LEPIWHI - *Lepidorhombus*
319 *whiffiagonis*) (Fig. 4). Overall best results were obtained by species sampled in the pelagic
320 survey, *Trisopterus esmarkii* (TRISESM) and *Merlangius merlangus* (MERNMER) were
321 exceptions to this, attaining also higher VEcv in bottom trawl survey. Inter-annual variability
322 of VEcv varied across species, with species (e.g. *Sardina pilchardus* - SARDPIL or *Trisopterus*
323 *minutus* - TRISMIN) showing very little change in the results among the years, whereas others
324 such as *Trisopterus luscus* (TRISLUS), showing a more variable results, although overall the
325 patterns observed in relation to each method, across the years were relatively stable, i.e. most
326 of the years within species showed a similar results (Fig. 5). Kriging based methods showed
327 very similar VEcv between each other, independently of the neighborhood considered (between
328 3 and 30 points, i.e. OK03-OK30), the fitting of the variogram manually (MKri) or
329 automatically (OKri), using depth as covariate (UKri) or even with conditional simulation
330 (CSim) (Fig. 5).
331 In 13% of the distributions the method with the highest VEcv showed no-overlapping of the
332 IQR with all remaining methods (23 cases representing 15 species, out of the 172 sp/year
333 distributions)(Fig. 5). Within those, in 10 cases, RFor was the best method, 6 cases it was
334 GAMd, 3 cases UKr, 2 cases IDW and 1 case TPS and MKr. *Sprattus sprattus* (SPRASPR) was
335 the only species from the pelagic survey with one distribution showing a clear winning method
336 on the first step. In none of the case studies all years for one species showed only one best
337 method. Thus, in the remaining 149 case studies a second selection step was required.

## SECOND SELECTION STEP: MULTIPLE CRITERIA

339 The best interpolation method according to each criteria varied widely across species-years,
340 confirming that different aspects of the distributions are taken into account by each measure
341 (Fig. 6; note that cases with ties were omitted). For example, RFor was the best method in terms
342 of a.pix_under, a.over.perc, CGdist, RMSE and MAE whereas UKri was the best according to
343 a.mean_perc and GAMd produced the maps with less deviations on the isotropy and with more
344 similar aggregation (Gini index, Fig. 6). It is also clear the contrast between the results given
345 by different kriging neighborhood in the best method by criteria (Fig. 6).
346 For a decision framework on the second step, all measures from each distribution were
347 integrated using a principal component analysis, as illustrated for four species/years in Fig. 7.

348 The variance explained by PC1 was always above 90%, although PC2 also proportionated
349 important information in discriminating the issues of the different interpolators relatively to the
350 measures considered. In the given examples, for *Argentina* sp. From 2014 (ARGENT.2014),
351 GAM and TPS methods showed highest deviations in terms of pixels under minimum whereas
352 kriging and Rfor showed highest spatial distortion, although integrating all criteria the best
353 method would be Kr30. Similar interpretation can be done for every species/year distribution,
354 and thus conclude that *Callionymus lyra* in 2012 (CLAMLYR.2012) best method would be
355 GAM, for *Conger conger* in 2015 (CONGCON.2015) would be Kr30 whereas for *Gadiculus*
356 *argenteus* in 2013 (GADIARG.2013) would be Kr7.
357 The results of the protocol on the interpolator selection procedure, given the two steps together
358 are found in Fig. 8. As already mentioned, few case studies were resolved on the first step (grey
359 boxes). The use of the multi-criteria privileged GAM, TPS, UKr and Rfor for the trawl data set
360 and TPS, OKri and UKr for the acoustic data set. It is also evident that the results differed
361 within methods (e.g. kriging neighborhood and type of kriging) showing that small
362 modifications on the methods can have a large impact on the surfaces produced. Further, it is
363 interesting to observe such a large disparity on the best method, not only between species but
364 also across years for the same species.

## Discussion

366 In the current work we develop a two-step procedure to aid researchers select the best
367 interpolation method for their data. The method uses a multi-criteria approach, that considers
368 error-based measures, changes in the spatial structure and data integrity after interpolation and
369 permits to determine in which particular aspect the interpolation is failing. The two-step
370 procedure was illustrated by comparing 20 interpolation methods applied to 175 distributions,
371 i.e. 35 species obtained during five years (2011-2015) of a typical bottom trawl survey and a
372 pelagic acoustic survey. In the first step of the selection procedure, all interpolation methods
373 within the highest VEcv's interquartile ranges are selected. In 13% of the case-studies no other
374 method had overlapping VEcv and thus, the selection is complete without having to go through
375 a second step. However, in the remaining 87% of the cases multiple methods were within the
376 best VEcv interquartile range and thus a second step was proposed, using additional criteria.
377 The effect of the interpolation was then integrated using PCA from which an index was
378 extracted to rank the quality of the interpolations.
379 In the current work we aimed to use the simplest and most available approach to summarise the indicators
380 in the second step, and this is why the PCA was selected. However, other multivariate statistical methods
381 besides PCA could be used as an alternative in future works. Furthermore, it is possible to use just the first
382 component (instead of two as in the current work) or to use the suggested indicators *per se*. Further work is
383 needed to develop this particular aspect.
384

385

## OTHER SELECTION PROCEDURES

387 Hengl *et al.* (2013) considered that the selection of a mapping procedure should account for
388 accuracy (considered to be measured by RMSE), bias (considered to be measured by MAE),
389 robustness (model sensitivity — in how many situations would the algorithm completely fail /
390 how much artifacts does it produces?), reliability (how good is the model in estimating the

391 prediction error, i.e. how accurate is the prediction variance considering the true mapping
392 accuracy?) and computation burden (the time needed to complete predictions). Visual
393 examination has been considered as equally important as accuracy measurements (Li et al.,
394 2011), although it is largely subjective and not explicitly defined, consistent or repeatable (Stow
395 et al., 2009). Model selection has involved evaluation, validation, performance, accuracy, skill,
396 efficiency or robustness, although these concepts are hard to disentangle and have been used
397 with different meanings on the literature (reviewed in Bellocchi et al., 2010). There is a clear
398 absence of a unifying selection procedure for interpolation models, although the most recent
399 works advocate the use of VEcv as a universal tool (Li, 2016, 2017). However, in 87% of the
400 distributions studied in the current work such approach was not sufficient to discriminate among
401 the interpolation methods, in view of the strong overlap between the respective distributions.
402 Further, the measures used for model selection should be explicit, with a straightforward
403 meaning and if possible, integrate multiple desirable properties of the interpolation procedure.
404 Similar to the current study, other authors have suggested that method selection should be multi-
405 criteria (Stow et al., 2009) and that the use of a single error measure may lead to incorrect
406 interpretation (Hoffman, 2015) .

## RMSE AND MAE

408 Each different measure comes with advantages and drawbacks. RMSE provides a measure of
409 error size, but it is sensitive to outliers as it places a lot of weight on large errors (Hernandez-
410 Stefanoni and Ponce-Hernandez, 2006). However, MAE and RMSE are among the best overall
411 measures of model performance as they summarize the mean difference in the units of observed
412 and predicted values (Willmott, 1982), although being highly correlated between each other,
413 with biomass/occurrence and algebraically related. Variance explained by predictive model
414 (VEcv) has been recently considered as the best error based criteria to evaluate interpolators
415 (Li, 2017, 2016). Our results indicate that this measure is straightforward to interpret and quick
416 evaluate thus it was included in the procedure, but for comparison with previous works, RMSE
417 and MAE were also incorporated. .

## CROSS VALIDATION

419 The error measures used to evaluate interpolation methods are traditionally estimated by cross
420 validation procedures, either leave one out (LOO)(Kilibarda et al., 2014) or k-fold cross-
421 validation (generally five or ten)(Davis, 1987; Li et al., 2011)(for a schematic overview of the
422 re-sampling strategies for model validation see Richter et al., 2012). First of all we decided to
423 not perform LOO cross-validation because in the case of skewness distributions and extreme
424 values of the input data, this kind of cross-validation might produce strange outputs (Hengl,
425 2009). Secondly in highly clustered spatial distribution (like the species distributions
426 considered), k (number of subsets, typically 5 or 10) should be large enough so that the data
427 into the k-subsets contains enough information on the whole model domain and the spatial
428 structure (Augustin et al., 2013). In our study, after a preliminary test, it was concluded that 10-
429 fold was a good compromise for marine species distribution (not reported for brevity).
430 Additionally, as the results of cross-validation strongly depended on the way the data is split
431 (folds), the process should be randomly repeated several times, and as consequence it is
432 obtained a distribution of the measure, which can then be used to compare the model's VEcv
433 (like it was done in the first step of the procedure). It is important to note also that cross-
434 validation is not necessarily independent, indeed, points used for cross-validation are subset of

435 the original sampling design. Consequently, if the original design is biased and/or non-
436 representative, then also the cross-validation might not reveal the true accuracy of a technique
437 (Hengl, 2009). Further, error-based measures estimated by cross-validation results can be
438 corrupted for clustered data sets on interpolator comparison (Hengl et al., 2013), highlighting
439 the importance of having additional criteria in the selection procedure when the decision is not
440 evident.
441 Cross validation can also be used to define the spatial model (Fortin and Dale, 2005; Gaetan et
442 al., 2010; Wackernagel, 1998) and kriging neighborhood (Paramo and Roa, 2003) within the
443 geostatistical methods. This is particularly crucial, because the effectiveness of the kriging
444 depends on how well the selected model fits the data (Fortin and Dale, 2005). There were only
445 small changes in the variance explained by the predictive model observed between kriging
446 computed with different neighborhood or to the process of defining the spatial model (manual
447 vs. automatic), when compared to other methods, but large changes were observed when the
448 other comparison criteria were included. For example, Gini index or the percentage of over
449 predictions changed widely with kriging neighborhood. Such comparison can also provide an
450 idea of the variation due to the parametrization between and within each different techniques,
451 as recommended in Araújo & Guisan (Araújo and Guisan, 2006)(within-model vs. between-
452 model comparisons). Thus, the proposed protocol can also be applied as a tool to improve model
453 specification by within model comparison, in future works. Similarly, the effect of a more
454 detailed parametrization of the other interpolation methods on the quality of the predictions
455 cannot be ignored. The method developed can also be used for such parametrization, as it was
456 explored in the current work for kriging.

457 SPATIAL AND DATA INTEGRITY

458 Spatial indicators have been develop with the aim of quantifying distribution's spatial patterns
459 (Bez and Rivoirard, 2001; Woillez et al., 2009b, 2007), but have also been applied as a model
460 validation tool, to compare the model's outputs with sampled data for example (e.g. Huret et
461 al., 2010)(Rufino et al., 2018). Additionally, these metrics are particularly sensitive to
462 interpolation (Rufino et al., 2019) and are well suited to assess the spatial integrity of the sample
463 data, after interpolation. The four metrics selected in our study highlight shifts in the main
464 spatial features of the distributions, namely its location (center of gravity), dispersion (inertia),
465 direction (isotropy) and aggregation (Gini index). It is expected (similar to what is done for
466 modelling procedures) that a better interpolator would cause a minimum effect on those spatial
467 aspects, when compared with the sampled data, therefore preserving the data spatial integrity.
468 Future works can assess the use of other spatial indicators, such as the index of collocation for
469 example, that may potentially be interesting with this aim.
470 Similarly, a good interpolation method should preserve the data limits of the sampled data. This
471 is often done within works of spatial analysis, but rarely mentioned and hardly quantified. All
472 those measures used in the selection procedure were made relative to the area and in the same
473 direction (i.e. the larger the value, the worst the mode) for comparability purposes.
474        These aspects together are of outmost importance for species distributions maps, and have
475 never been systematized previously. It can be argued that some of these measures just report
476 the intrinsic properties of the interpolation methods and thus could be inferred solely on
477 theoretical grounds. For example, kriging methods tend to under-estimate the maximum
478 biomass (Bivand et al., 2013; Cressie, 1993). The proposed measures evidence those properties,
479 and make them accessible without requiring a strong expertise on spatial analyses. On the other
480 way, some measures, for example those with reference to the spatial integrity, are much less
481 evident to describe theoretically.

482    PURELY GEOGRAPHIC METHODS *VS*. METHODS USING COVARIATES

483    Purely geographic methods, i.e. using only geographic coordinates to produce spatial
484    predictions (i.e. ordinary kriging, TPS, IDW, etc.) are essential for the cases where there is a
485    belief that geographic processes are dominant over environmental ones or in the absence of
486    adequate environmental predictors (Elith and Leathwick, 2009). In the majority of cases the
487    purpose of the statistical modelling is the prediction of species distribution, whereas the
488    relationships between species and the environment tend to be a secondary consideration
489    (Austin, 2002; Guisan and Zimmermann, 2000). This is also the focus of the current work and
490    probably the commonest situation in marine studies or fisheries management. Thus, the applied
491    models with covariates used only topographic variables directly derivable from depth, which is
492    widely available. Unlike on its terrestrial counterpart, the effect of bottom topography on fish
493    distribution is seldom tackled in fisheries ecology (Giannoulaki et al., 2006, 2003). These
494    features are further advantageous for being relatively stable through time on these areas
495    (Maravelias, 1999)(unlike other environmental characteristics) thus being potentially
496    interesting also for long term studies, where other variables are not available.
497    It can be expected that topographic information of the sea bottom is more important for
498    demersal species than for the pelagic ones. However, sea bottom topography features are known
499    to be determinant for small pelagics species (Giannoulaki et al., 2006; Maravelias, 1999) and
500    some of these species occur also near the sea bottom (e.g. *Scomber* sp. and *Trachurus*
501    *trachurus*). However, the models done with topographic variables were not better than those
502    using just geographic variables for any pelagic species, but it is imperative to see the marine
503    environment as a continuous system where all aspects are connnected, and thereby only
504    manageable through an ecosystem approach (Cotter et al., 2009; Doray et al., 2018). Other
505    relevant environmental variables such as temperature and productivity, would improve the
506    models with covariates, but can be more species specific. When a set of ecological covariates
507    is available, whatever these are, the current method is also applicable.
508

509    SPATIAL AUTOCORRELATION

510    The notion of spatial autocorrelation is largely attributed to Tobler's 1st Law of Geography,
511    "Everything is related to everything else, but near things are more related than distant things"
512    (Tobler, 1970). Spatial autocorrelation is widely present in marine species species distribution
513    and is an essential aspect to acount for in spatial prediction (Dormann et al., 2007; Elith and
514    Leathwick, 2009; Legendre, 1993). We verified using a large empirical data set that models
515    accounting for spatial autocorrelation (i.e. geostatitical models) showed higher VEcv overall
516    both for the bottom trawl data and for the acoustic pelagic survey. This difference was more
517    pronounced on the acoustic pelagic data, where the number of samples is also much higher and
518    the spatial models of the variogram, better defined (pers. obs. MMR). Improving the spatial
519    model definition increases the effectiveness of kriging (Fortin and Dale, 2005). This was also
520    observed in Rufino *et al* (2006) using simulated data, where the precision and accuracy of the
521    kriging predictions increased with the sample size, as well as the importance of spatial
522    autocorrelation. On the other way, in some cases, high data variability may hamper the retrieval
523    of the spatial models and mask spatial autocorrelation (Rufino et al., 2006).
524    The clumped spatial patterns typical of marine species distribution can emerge simply as a result
525    of the spatial autocorrelation of the environmental and of biotic processes (Legendre, 1993). In
526    the current work, most species evidenced the presence of auto-correlation in the experimental

variogram model. Strong residual geographic patterning generally indicates that either key environmental predictors are missing, that the model is mis-specified or that geographic factors are influential (Elith and Leathwick, 2009). In a study as broad as the current one this would be a natural consequence as it was not the aim to explore the key environmental predictors of each species, nether to parametrise in detail each method. Recent works have shown excellent results on the application of combined methods (random forest + kriging) in other areas (Appelhans et al., 2015; Diesing et al., 2014; Hengl et al., 2015; Li et al., 2011, 2013, 2016), and thus it would be interesting to explore such applications to marine SDM on future works. Further, the selection protocol is extensible interpolation methods on the spatio-temporal domain.

It is evident that the best interpolation changed widely across species and years, and thus in each case a detailed analysis is required. Furthermore, the fact that certain interpolation models performed better for some species in a certain dataset, does not imply that it will always perform better with other fisheries datasets (Davis, 1987). Nevertheless, our application of the selection protocol to the two surveys reveals general guidelines for the variability of the results given by different interpolation methods. It is clear that each model need to be parameterized in detail, individually and according to the species data for a proper spatial analysis and that neglecting ecological knowledge is a limiting factor in the use of statistical modelling to predict species distribution (Austin, 2002).

We conclude that the proposed 2-step approach for method's selection has several benefits: 1) it is accessible and does not require any complex software or sophisticated method; 2) it is explicit in the sense that it evidences the benefits of each interpolation model relative to the others, empirically, that is on the maps produced and thus, permits to make a decision accordingly to the objectives of each study; 3) it takes into account different criteria, thus integrating several desirable properties of interpolation methods; 4) it does not depend on the method used or the data characteristics, thus being universal and can be applied to virtually any method developed in the future.

## Acknowledgments

## Data availability

The data used in the current work can be located in

https://campagnes.flotteoceanographique.fr/series/8/

## References

Aalto, J., Pirinen, P., Heikkinen, J., Venäläinen, A., 2013. Spatial interpolation of monthly climate data for Finland: comparing the performance of kriging and generalized additive models. Theor. Appl. Climatol. 112, 99–111. doi:10.1007/s00704-012-0716-9

Amante, C.J., Eakins, B.W., 2016. Accuracy of interpolated bathymetry in digital elevation models. J. Coast. Res. 76, 123–133. doi:10.2112/SI76-011

Appelhans, T., Mwangomo, E., Hardy, D.R., Hemp, A., Nauss, T., 2015. Evaluating machine learning approaches for the interpolation of monthly air temperature at. Spat. Stat. 14, 91–113. doi:10.1016/j.spasta.2015.05.008

Araújo, M.B., Guisan, A., 2006. Five (or so) challenges for species distribution modelling. J. Biogeogr. 33, 1677–1688. doi:10.1111/j.1365-2699.2006.01584.x

Araujo, M.B., Peterson, A.T., 2012. Uses and misuses of bioclimatic envelope modeling. Ecology 93, 1527–1539. doi:10.1890/11-1930.1

Augustin, N.H., Trenkel, V.M., Wood, S.N., Lorance, P., 2013. Space-time modelling of blue ling for fisheries stock management. Environmetrics 24, 109–119. doi:10.1002/env.2196

Austin, M., 2007. Species distribution models and ecological theory: A critical assessment and some possible new approaches. Ecol. Modell. 200, 1–19. doi:10.1016/j.ecolmodel.2006.07.005

Austin, M.P., 2002. Spatial prediction of species distribution: an interface between ecological theory and statistical modelling. Ecol. Modell. 157, 101–118. doi:10.1016/S0304-3800(02)00205-3

Baddeley, A., Gregori, P., Mateu, J., Stoica, R., Stoyan, D., Bickel, E.P., Diggle, P., Fienberg, S., Gather, U., Baddeley, A., Stoyan, D., 2006. Case Studies in Spatial Point Process Modeling. Springer. doi:10.1007/0-387-31144-0

Bellocchi, G., Rivington, M., Donatelli, M., Matthews, K., 2010. Validation of biophysical models : issues and methodologies. Agron. Sustain. Dev. 30, 109–130. doi:10.1051/agro/2009001

Bez, N., Rivoirard, J., 2001. Transitive geostatistics to characterise spatial aggregations with diffuse limits: an application on mackerel ichtyoplankton. Fish. Res. 50, 41–58. doi:10.1016/S0165-7836(00)00241-1

Bivand, R., Pebesma, E., Gómez-Rubio, V., 2013. Applied spatial data analysis with R, Use R. Springer. doi:10.1007/978-0-387-78171-6

Cotter, J., Petitgas, P., Abella, A., Apostolaki, P., Mesnil, B., Politou, C.-Y., Rivoirard, J., Rochet, M.-J., Spedicato, M.T., Trenkel, V.M., Woillez, M., 2009. Towards an ecosystem approach to fisheries management (EAFM) when trawl surveys provide the main source of information. Aquat. Living Resour. 22, 243–254. doi:10.1051/alr/2009025

Cressie, N.A.C., 1993. Statistics for spatial data, Statistics for Spatial Data. John Wiley &

Sons, New York. doi:10.1002/9781119115151

Dauvin, J.-C., Thiebaut, E., Gesteira, J.L.G., Ghertsos, K., Gentil, F., Ropert, M., Sylvand, B., 2004. Spatial structure of a subtidal macrobenthic community in the Bay of Veys (western Bay of Seine, English Channel). J. Exp. Mar. Bio. Ecol. 307, 217–235. doi:10.1016/j.jembe.2004.02.005

Davis, B.M., 1987. Uses and abuses of cross-validation in geostatistics. Math. Geol. 19, 241–248. doi:10.1007/BF00897749

Diesing, M., Green, S.L., Stephens, D., Lark, R.M., Stewart, H.A., Dove, D., 2014. Mapping seabed sediments: Comparison of manual, geostatistical, object-based image analysis and machine learning approaches. Cont. Shelf Res. 84, 107–119. doi:10.1016/j.csr.2014.05.004

Doray, M., Duhamel E. , Huret M. , Petitgas P., M.J., 2002. PELGAS. doi:10.18142/18

Doray M., Badts V., Masse J., Duhamel E., Huret M., Doremus G., P.P., 2014. Manuel des protocoles de campagne halieutique. Campagnes PELGAS (PELagiques GAScogne) / Manual of fisheries survey protocols. PELGAS surveys (PELagiques GAScogne). doi:10.13155/30259

Doray, M., Petitgas, P., Romagnan, J.B., Huret, M., Duhamel, E., Dupuy, C., Spitz, J., Authier, M., Sanchez, F., Berger, L., Dorémus, G., Bourriau, P., Grellier, P., Massé, J., 2018. The PELGAS survey: Ship-based integrated monitoring of the Bay of Biscay pelagic ecosystem. Prog. Oceanogr. 166, 15–29. doi:10.1016/j.pocean.2017.09.015

Dormann, C.F., McPherson, J.M., Araújo, M.B., Bivand, R., Bolliger, J., Carl, G., Davies, R.G., Hirzel, A., Jetz, W., Kissling, W.D., Kuhn, I., Ohlemuller, R., Peres-Neto, P.R., Reineking, B., Schro¨der, B., Schurr, F.M., Wilson, R., 2007. Methods to account for spatial autocorrelation in the analysis of species distributional data: a review. Ecography (Cop.). doi: 10.11. doi:10.1111/j.2007.0906-7590.05171.x

Elith, J., Leathwick, J.R., 2009. Species distribution models: ecological explanation and prediction across space and time. Annu. Rev. Ecol. Evol. Syst. 40, 677–697. doi:10.1146/annurev.ecolsys.110308.120159

Evaluation Halieutique de l'Ouest Européen, EVHOE cruise, RV Thalassa, IFREMER, n.d. doi:10.18142/8

Evans JS, 2017. spatialEco. R package version 0.0.1-7.

Fortin, M.M.-J., Dale, M.R.T., 2005. Spatial analysis: a guide for ecologists, 4th ed. Cambridge University Press, Cambridde. doi:10.2277/0521804345

Gaetan, C., Guyon, X., Bleakley, K., 2010. Spatial Statistics and Modeling, Media. Springer. doi:10.1007/978-0-387-92257-7

Gasch, C.K., Hengl, T., Gräler, B., Meyer, H., Magney, T.S., Brown, D.J., 2015. Spatio-temporal interpolation of soil water, temperature, and electrical conductivity in 3D + T: The Cook Agronomy Farm data set. Spat. Stat. 14, 70–90. doi:10.1016/j.spasta.2015.04.001

Giannoulaki, M., Machias, A., Koutsikopoulos, C., Haralabous, J., Somarakis, S., Tsimenides, N., 2003. Effect of coastal topography on the spatial structure of the populations of small pelagic fish. Mar. Ecol. Prog. Ser. 265, 243–253. doi:10.3354/meps265243

Giannoulaki, M., Machias, A., Koutsikopoulos, C., Somarakis, S., 2006. The effect of coastal topography on the spatial structure of anchovy and sardine. ICES J. Mar. Sci. 63, 650–662. doi:10.1016/j.icesjms.2005.10.017

Guisan, A., Edwards Thomas C., J., Hastie, T., 2002. Generalized linear and generalized additive models in studies of species distributions: setting the scene. Ecol. Modell. 157, 89–100. doi:10.1016/S0304-3800(02)00204-1

Guisan, A., Zimmermann, N.E., 2000. Predictive habitat distribution models in ecology. Ecol. Modell. 135, 147–186. doi:10.1016/S0304-3800(00)00354-9

652    Hengl, T., 2009. A Practical guide to Geostatistical Mapping, Scientific and Technical
653        Research series. doi:10.1016/0277-9390(86)90082-8
654    Hengl, T., Geuvelink, G.B.M., Stein, A., 2003. Comparison of kriging with external drift and
655        regression kriging. Technical note, ITC, Available on-line at
656        http://www.itc.nl/library/Academic_output/. doi:10.1016/S0016-7061(00)00042-2
657    Hengl, T., Heuvelink, G., Stien, A., 2004. A generic frame work for the spatial prediction of
658        soil variables based on regression-kriging. Geoderma 120, 75–93.
659        doi:10.1016/j.geoderma.2003.08.018
660    Hengl, T., Heuvelink, G.B.M., Kempen, B., Leenaars, J.G.B., Walsh, M.G., Shepherd, K.D.,
661        Sila, A., MacMillan, R.A., De Jesus, J.M., Tamene, L., Tondoh, J.E., 2015. Mapping soil
662        properties of Africa at 250 m resolution: Random forests significantly improve current
663        predictions. PLoS One 10, 1–26. doi:10.1371/journal.pone.0125814
664    Hengl, T., Heuvelink, G.B.M., Rossiter, D.G., 2007. About regression-kriging: From
665        equations to case studies. Comput. Geosci. 33, 1301–1315.
666        doi:10.1016/j.cageo.2007.05.001
667    Hengl, T., MacMillan, R.A., Nikolić, M., 2013. Mapping efficiency and information content.
668        Int. J. Appl. Earth Obs. Geoinf. 22, 127–138. doi:10.1016/j.jag.2012.02.005
669    Hoffman, S., 2015. Assessment of prediction accuracy in autonomous air quality models.
670        Desalin. Water Treat. 57, 1322–1326. doi:10.1080/19443994.2014.1002283
671    Hui, F.K.C., Warton, D.I., Foster, S.D., Dunstan, P.K., 2013. Species Distribution Models:
672        Ecological Explanation and Prediction Across Space and Time vs. separate species
673        distribution models. Ecology 94, 1913–1919. doi:10.1890/12-1322.1
674    Huret, M., Petitgas, P., Woillez, M., 2010. Dispersal kernels and their drivers captured with a
675        hydrodynamic model and spatial indices: A case study on anchovy (Engraulis
676        encrasicolus) early life stages in the Bay of Biscay. Prog. Oceanogr. 87, 6–17.
677        doi:10.1016/j.pocean.2010.09.023
678    Kilibarda, M., Hengl, T., Heuvelink, G.B.M., Gräler, B., Pebesma, E., Perčec Tadič, M.,
679        Bajat, B., 2014. Spatio-temporal interpolation of daily temperatures for global land areas
680        at 1 km resolution. J. Geophys. Res. Atmos. 119, 2294–2313.
681        doi:10.1002/2013JD020803
682    Lark, M., Dove, D., Green, S., Stewart, H., Marchant, B., Diesing, M., 2016. Uncertainty in
683        predictions of seabed sediment classes based on grab samples and acoustic data.
684        Geophys. Res. Abstr. 18.
685    Lecours, V., Dolan, M.F.J., Micallef, A., Lucieer, V.L., 2016. A review of marine
686        geomorphometry, the quantitative study of the seafloor. Hydrol. Earth Syst. Sci. 20,
687        3207–3244. doi:10.5194/hess-20-3207-2016
688    Legates, D.R., McCabe Jr., G.J., 1999. Evaluating the use of "goodness of fit" measures in
689        hydrologic and hydroclimatic model validation. Water Resour. Res. 35, 233–241.
690        doi:10.1029/1998WR900018
691    Legendre, P., 1993. Spatial autocorrelation: trouble or new paradigm? Ecology 74, 1659–
692        1673. doi:10.2307/1939924
693    Legendre, P., Legendre, L., 1998. Numerical ecology, 2nd ed. Elsevier, Amsterdam.
694        doi:10.1016/S0304-3800(00)00291-X
695    Li, J., 2017. Assessing the accuracy of predictive models for numerical data: Not r nor r 2 ,
696        why not? Then what? PLoS One 12, 1–16. doi:10.1371/journal.pone.0183250
697    Li, J., 2016. Assessing spatial predictive models in the environmental sciences: Accuracy
698        measures, data variation and variance explained. Environ. Model. Softw. 80, 1–8.
699        doi:10.1016/j.envsoft.2016.02.004
700    Li, J., Heap, A.D., 2011. A review of comparative studies of spatial interpolation methods in
701        environmental sciences: Performance and impact factors. Ecol. Inform. 6, 228–241.

702        doi:10.1016/j.ecoinf.2010.12.003

703 Li, J., Heap, A.D., 2008. A review of spatial interpolation methods for environmental
704        scientists, Geoscience Australia.

705 Li, J., Heap, A.D., Potter, A., Huang, Z., Daniell, J.J., 2011. Can we improve the spatial
706        predictions of seabed sediments? A case study of spatial interpolation of mud content
707        across the southwest Australian margin. Cont. Shelf Res. 31, 1365–1376.
708        doi:10.1016/j.csr.2011.05.015

709 Li, J., Justy, P., Siwabessy, W., Tran, M., Huang, Z., Heap, A.D., 2013. Predicting Seabed
710        Hardness Using Random Forest in R, in: Elsevier (Ed.), Data Mining Applications with
711        R. pp. 299–329. doi:10.1016/B978-0-12-411511-8.00011-6

712 Li, J., Siwabessy, J., Tran, M., 2016. Selecting optimal random forest predictive models : a
713        case study on predicting the spatial distribution of seabed hardness selecting optimal
714        random forest predictive models: a case study on predicting the spatial distribution of
715        seabed hardness. PLoS One 11, 1–29. doi:10.1371/journal.pone.0149089

716 Maravelias, C.D., 1999. Habitat selection and clustering of a pelagic fish: effects of
717        topography and bathymetry on species dynamics. Can. J. Fish. Aquat. Sci. 56, 437–450.
718        doi:10.1139/f98-176

719 Morfin, M., Bez, N., Fromentin, J.M., 2016. Habitats of ten demersal species in the Gulf of
720        Lions and potential implications for spatial management. Mar. Ecol. Prog. Ser. 547, 219–
721        232. doi:10.3354/meps11603

722 Naimi, B., Araújo, M.B., 2016. Sdm: A reproducible and extensible R platform for species
723        distribution modelling. Ecography (Cop.). 39, 368–375. doi:10.1111/ecog.01881

724 Nash, J.E., Sutcliffe, J. V, 1970. River flow forecasting through conceptual models part I-a
725        discussion of principles. J. Hydrol. 10, 282–290. doi:10.1016/0022-1694(70)90255-6

726 Nychka, D., 2016. fields: Tools for Spatial Data. doi:10.5065/D6W957CT

727 Olden, J.D., Jackson, D.A., 2002. A comparison of statistical approaches for modelling fish
728        species distributions. Freshw. Biol. 47, 1976–1995. doi:10.1046/j.1365-
729        2427.2002.00945.x

730 Paramo, J., Roa, R., 2003. Acoustic-geostatistical assessment and habitat-abundance relations
731        of small pelagic fish from the Colombian Caribbean. Fish. Res. 60, 309–319.
732        doi:10.1016/S0165-7836(02)00142-X

733 Richter, K., Atzberger, C., Hank, T.B., Mauser, W., 2012. Derivation of biophysical variables
734        from Earth observation data: validation and statistical measures. J. Appl. Remote Sens.
735        6, 063557–1. doi:10.1117/1.JRS.6.063557

736 Rufino, M.M., Bez, N., Brind'Amour, A., 2019. Influence of data pre-processing on the
737        behavior of spatial indicators. Ecol. Indic. 99, 108–117.
738        doi:10.1016/j.ecolind.2018.11.058

739 Rufino, M.M., Bez, N., Brind'Amour, A., 2018. Integrating spatial indicators in the
740        surveillance of exploited marine ecosystems. PLoS One 13, e0207538.
741        doi:10.1371/journal.pone.0207538

742 Rufino, M.M., Stelzenmüller, V., Maynou, F., Zauke, G.P., 2006. Assessing the performance
743        of linear geostatistical tools applied to artificial fisheries data. Fish. Res. 82, 263–279.
744        doi:10.1016/j.fishres.2006.06.013

745 Sluiter, R., 2009. Interpolation methods for climate data literature review, KNMI intern
746        rapport.

747 Stow, C.A., Jolliff, J., McGillicuddy, D.J., Jr.,  c S.C.D., Allen, J.I., Friedrichs, M.A.M.,
748        Rose, K.A., Wallhead, P., 2009. Skill assessment for coupled biological/physical models
749        of marine systems. J. Mar. Syst. 76, 4–15. doi:10.1016/j.jmarsys.2008.03.011

750 Thorson, J.T., Shelton, A.O., Ward, E.J., Skaug, H.J., 2015. Geostatistical delta-generalized
751        linear mixed models improve precision for estimated abundance indices for West Coast

752          groundfishes. ICES J. Mar. Sci. 72, 1297–1310. doi:10.1093/icesjms/fsu243

753    Tobler, W.R., 1970. A computer movie simulation urban growth in Detroit region. Econ.
754          Geogr. 46, 234–240. doi:10.1126/science.11.277.620

755    Wackernagel, H., 1998. Multivariate geostatistics: An introduction with applications., Second.
756          ed. Springer-Verlag, Berlin. doi:10.1007/978-3-662-03550-4

757    Webster, R., Oliver, M., 2007. Geostatistics for environmental scientists, 2nd ed. ed. John
758          Wiley and Sons, New York. doi:10.1002/9780470517277

759    Willmott, C.J., 1982. Some comments on the evaluation of model performance. Bull. Am.
760          Meteorol. Soc. 1309–1313. doi:10.1175/1520-0477(1982)063<1309:SCOTEO>2.0.CO;2

761    Willmott, C.J., 1981. On the validation of models. Phys. Geogr. 2, 184–194.
762          doi:10.1080/02723646.1981.10642213

763    Willmott, C.J., Robeson, S.M., Matsuura, K., 2012. A refined index of model performance.
764          Int. J. Climatol. 32, 2088–2094. doi:10.1002/joc.2419

765    Willmott, C.J., Robeson, S.M., Matsuura, K., Ficklin, D.L., 2015. Assessment of three
766          dimensionless measures of model performance. Environ. Model. Softw. 73, 167–174.
767          doi:10.1016/j.envsoft.2015.08.012

768    Wilson, M.F.J., Connell, B.O., Guinan, J.C., Grehan, A.J., 2007. Multiscale terrain analysis of
769          multibeam bathymetry data for habitat mapping on the continental slope.
770          doi:10.1080/01490410701295962

771    Woillez, M., Poulard, J.-C., Rivoirard, J., Petitgas, P., Bez, N., 2007. Indices for capturing
772          spatial patterns and their evolution in time, with application to European hake
773          (*Merluccius merluccius*) in the Bay of Biscay. ICES J. Mar. Sci. 64, 537–550.
774          doi:10.1093/icesjms/fsm025

775    Woillez, M., Rivoirard, J., Fernandes, P.G., 2009a. Evaluating the uncertainty of abundance
776          estimates from acoustic surveys using geostatistical simulations. ICES J. Mar. Sci. 66,
777          1377–1383. doi:10.1093/icesjms/fsp137

778    Woillez, M., Rivoirard, J., Petitgas, P., 2009b. Notes on survey-based spatial indicators for
779          monitoring fish populations. Aquat. Living Resour. 22, 155–164.
780          doi:10.1051/alr/2009017

781    Wood, N.S., 2006. Generalized Additive Models: An Introduction with R, Texts in Statistical
782          Science Series. Chapman & Hall/CRC. doi:10.1201/9781315370279

783    Zuur, A.F., Ieno, E.N., Smith, G.M., 2007. Analysing Ecological Data, Profiles of drug
784          substances, excipients, and related methodology. Springer. doi:10.1016/B978-0-12-
785          387667-6.00013-0

786

787

FIGURE LEGENDS:

789    Fig. 1: Graphical abstract: Conceptual diagram of the method developed to compare the
790    interpolators.
791    Fig. 2: Frequency of the variance explained by predictive model classes between spatial
792    interpolation methods, for the bottom trawl survey (EVHOE, upper panel) and the acoustic
793    pelagic survey (PELGAS, lower panel). Find further details on the methods codes in text.
794    Fig. 3 : Variance explained by the predictive model between spatial interpolation method's
795    main families, by survey (blue triangles for pelagic survey, PELGAS and red bals for bottom
796    trawl survey, EVHOE)(mean and respective 95% CI estimated using bootstrap). LM: linear
797    model; GAM: generalised additive models; IDW: inverse distance weighting; Vor: voronoi
798    triangulation; TPS: thin plate spline; Kri: kriging and conditional simulation; Covar: multiple
799    regression, regretion tree and random forest (simple and mixed, i.e. with kriged residuals).
800    Please find further details on the methods codes in the text.
801    Fig. 4 : Variance explained by the predictive model between species, for interpolation methods
802    using topographic covariates (orange line with squares, IMCov) and for methods just using
803    geographic coordinates and eventually depth (green line with diamonds, IMGeo). Filled
804    symbols represent the species captured in the bottom trawl survey whereas open symbols
805    indicate the acoustic pelagic survey. Mean and respective 95% CI estimated using bootstrap is
806    represented. Species were ordered by IMGeo VEcv. Find further details on the species codes
807    in the text.
808    Fig. 5: Variance explained by the model of each interpolation method (median), estimated by
809    cross validation for all species-year distributions. Red points indicate that the interpolation
810    method was within the best method interquartile range and red star indicate the best VEcv in
811    each case whereas black dots indicate models that did not passed for thr second step. Grey
812    shadded area correspond to methods carried out using the 11 topographic covariates (IMCov),
813    whereas white background shows methods using only lat+long and some depth (IMGeo). Please
814    find further details on the methods and species codes in the text.
815    Fig. 6: Distributions and methods that required the second step. Winning interpolation method
816    according to each measure criteria (left panel; only cases where more than one method showing
817    its VEcv within the highest method VEcv interquartile range were selected and situations with
818    ties were excluded, i.e. several methods showing the same classification according to the
819    criteria). Winning method for each measure-criteria by species-year distribution (right panel).
820    Vertical grey line separates the methods using several covariates from the others. Please find
821    further details on the criteria and methods codes in the text. The colour legend is represented in
822    the barplot.
823    Fig. 7: Second step of the spatial interpolator's selection protocol applied to four case studies
824    (2 bottom trawl, EVHOE and 2 acoustic pelagic, PELGAS). On the right side plots, the PCA
825    shows where each interpolation method (represented with circles, orange-red coloured,
826    according to the distance to the centre) failed according to the measures representing the three
827    selection criteria (error-based in violet, spatial integrity in green and data integrity in blue). On
828    the left panels, the inverse Euclidean distance to the centre of each method, provides the
829    quantitative decision integrative measure. Please find the details of the code's labels in the text.
830    Fig. 8: Winning spatial interpolation method among the different approaches considered, for
831    each case study (species-year distributions) according to the two step selection procedure (1st
832    step using IQR VEcv identified with orange line and 2nd step, using the 3 criteria with 11
833    measures, identified with blue line). Number of cases of each selected method by survey.
834    Vertical grey line separates the methods using several covariates from the others. See further
835    details of the species codes and methods on text. The colour legend is represented in the barplot.

836
837

838 **Figures**

839 **Fig. 1**
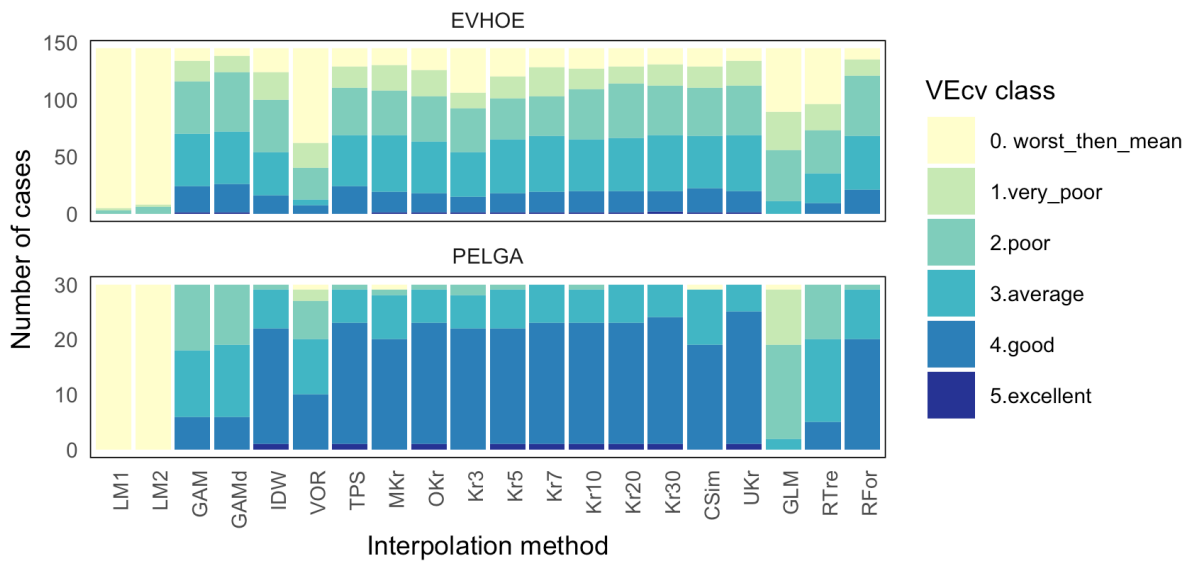


840
841

842 **Fig. 2: VEcv classes + RT**



843
844
845

846 **Fig. 3: VEcv by family**



847

848 **Fig. 4: VEcv by Sp**

849
850

851

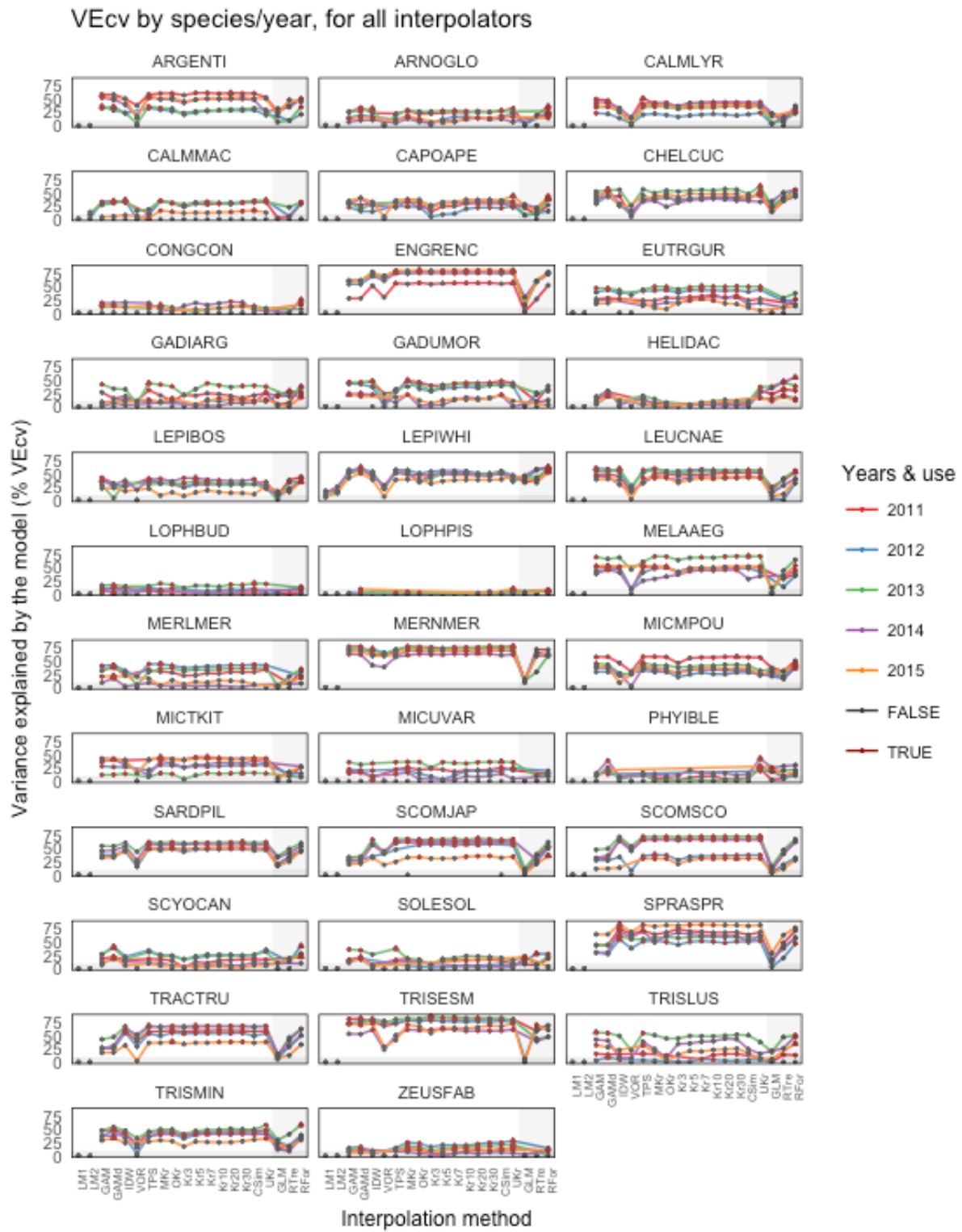**Fig. 5: all VEcvs**

VEcv by species/year, for all interpolators

856 **Fig. 6: Best method according to each criteria (2nd step only).**



857

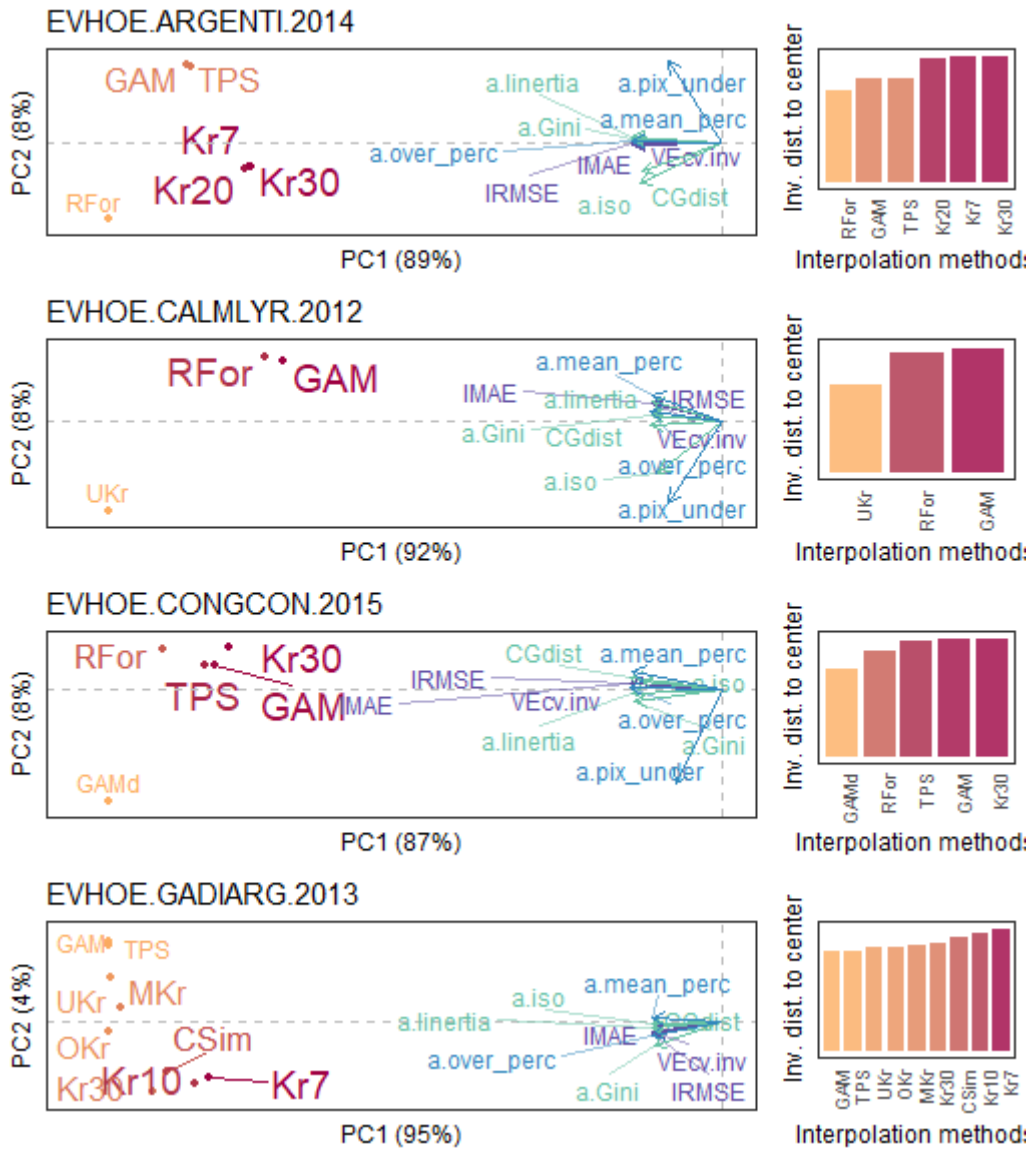858     **Fig. 7: second step PCA examples**



859
860

861     Fig. 8: **Winning methods overall;**

867     **Supplement 1**

868

869     Location of the sampling stations from 2015 of the bottom trawl survey (EVHOE, left panel),
870     of the pelagic acoustic survey (PELGAS) on the Gulf of Gascogne/Bay of Biscay in the North
871     Atlantic coast of France (right panel). Number of hauls per year on EVHOE, from 2011 to 2015
872     were respectively 220, 195, 208, 219 and 148.



873
874     **Species list names and abreviations**
875     EVHOE: *Argentina* sp. (ARGENTI), *Arnoglossus* sp. (ARNOGLO), *Callionymus lyra*
876     (CALMLYR), *Callionymus maculatus* (CALMMAC), *Capros aper* (CAPOAPE),
877     *Chelidonichthys cuculus* (CHELCUC), *Conger conger* (CONGCON), *Eutrigla gurnardus*
878     (EUTRGUR), *Gadiculus argenteus* (GADIARG), *Gadus morhua* (GADUMOR), *Helicolenus*
879     *dactylopterus* (HELIDAC), *Lepidorhombus boscii* (LEPIBOS), *Lepidorhombus whiffiagonis*
880     (LEPIWHI), *Leucoraja naevus* (LEUCNAE), *Lophius budegassa* (LOPHBUD), *Lophius*
881     *piscatorius* (LOPHPIS), *Melanogrammus aeglefinus* (MELAAEG), *Merluccius merluccius*
882     (MERLMER), *Merlangius merlangus* (MERNMER), *Microchirus variegatus* (MICUVAR),
883     *Micromesistius poutassou* (MICPOU), *Microstomus kitt* (MICKIT), *Phycis blennoides*
884     (PHYIBLE), *Scyliorhinus canicula* (SCYOCAN), *Solea solea* (SOLESOL), *Trisopterus*
885     *esmarkii* (TRISESM), *Trisopterus luscus* (TRISLUS), *Trisopterus minutus* (TRISMIN) and
886     *Zeus faber* (ZEUSFAB)
887     PELGAS: *Engraulis encrasicolus* (ENGRENC), *Sardina pilchardus* (SARDPIL), *Scomber*
888     *japonicas* (SCOMJAP), *Scomber scombrus* (SCOMSCO), *Sprattus sprattus* (SPRASSPR) and
889     *Trachurus trachurus* (TRACTRU)

890

891     **Details of the survey**:
892     The EVHOE survey has been carried out on the R/V Thalassa, a stern trawler of 73.7 m length
893     by 14.9 m wide (tonnage of 3022 t). The fishing gear used is a GOV 36/47 without exocet Kite
894     which is replaced by 6 additional floats and with a horizontal opening of 20 m and a vertical
895     opening of 4 m.


896     **Supplement 2**

897

898     **Description of terrain variables:**
899     Slope represent the terrain steepness (arrangement and magnitude of elevation
900     differences)(slope) whereas terrain aspect (aspect) measures its orientation in degrees, relative
901     to the north and it is particularly important to exposure to currents or water movement (Wilson

902 et al., 2007). From aspect, northerness and easterness were derived (Wilson et al., 2007). Profile
903 curvature defines convex/concave areas, represented by the rate of slope change along a profile,
904 i.e. the surface of the steepest down-slope direction (surf.curv)(package 'spatialEco')(Evans JS,
905 2017). Bathymetric Position Index is the difference between the value of a cell and the mean
906 value of its surrounding cells and provides an indication of whether any particular pixel forms
907 part of a positive (e.g., crest) or negative a (e.g., trough) feature of the surrounding terrain (TPI,
908 marine version of the topographic position index). Terrain Ruggedness Index represents terrain
909 variability whereas roughness represents the bathymetric amplitude of a cell and its
910 surroundings (TRI). Surface flow confluence indicates the steepest downhill path (flowdir).
911 Local Moran I was calculated as a measure of local spatial autocorrelation in the bathymetric
912 neighborhood (moran)(Diesing et al., 2014; Li et al., 2016). Additionally, distance to the nearest
913 coast was also estimated (dist.coast).
914

915 **Maps of the topographic variables used in the current work (bottom trawl area, EVHOE)**
916 After gridding all variables to EVHOE prediction grid (1117 pixels, 15 km).

917



918 **Supplement 3**

919

920 **Table 1: Summary of the interpolation methods used in the current work.**

| Interpolator | Description | R-package /Function | Covar |
|---|---|---|---|
| IMGeo | | | |
| LM1 | 1st-order trend surface | gstat::gstat(vari ~ 1, degree=1) | - |
| LM2 | 2nd order trend surface | gstat::gstat(vari ~ 1, degree=2) | - |
| **GAM** | Generalized Additive Model | | |

| | | | |
|---|---|---|---|
| | (Wood, 2006) | | |
| GAM | GAM in function of lat and long (in UTM) | mgcv ::gam (lvari ~ s(lat ,long)) | |
| GAMd | GAM in function of lat and long (in UTM) and depth | mgcv ::gam (lvari ~ s(lat ,long)+s(Depth)) | Depth |
| **IDW** | Inverse distance weight | | |
| IDW | Optimized using cross validation | gstat::gstat(lvari~1, nmax=opt$par[1], set=list(idp=opt$par[2])) | - |
| **VOR** | Voronoi thesselation (Fortin and Dale, 2005) | dismo::voronoi(dat.s["lvari"]) | |
| **TPS** | Thin Plate Spline interpolation (Nychka, 2016) | | |
| **Kriging** | Methods using kriging (Bivand et al., 2013) | | |
| MKr | Ordinary kriging interpolation (manual fitting of variogram) | gstat::krige(lvari ~ 1) | - |
| CSim | Stochastic conditional simulation | gstat::krige(lvari ~ 1, nsim = 1000, nmax = 20) | - |
| OKri | | automap::autoKrige(lvari~1, model = c("Sph", "Exp")))(gstat) | - |
| 17-22. Kr3 Kr5 Kr7 Kr10 Kr20 Kr30 | OK with various neighborhood | automap::autoKrige(lvari~1, model = c("Sph", "Exp"), nmax=jj) (gstat) | - |
| UKr | Universal kriging, with depth as covariate | automap::autoKrige(lvari~Depth, model = c("Sph", "Exp")) (gstat) | Depth |
| GLM | Multiple regression | dat.s, lvari~Depth+slope+aspect+eastness+northness+surf.curv+TPI+TRI+rough+dist.coast+flowdir | Depth+slope+aspect+eastness+northness+surf.curv+TPI+TRI+rough+dist.coast+flowdir |
| RTre | Regression tree | rpart:: rpart ( dat.s, vari ~ Depth+slope+aspect+eastness+northness+surf.curv+TPI+TRI+rough+dist.coast+flowdir) | Depth+slope+aspect+eastness+northness+surf.curv+TPI+TRI+rough+dist.coast+flowdir |
| RFor | Random forest | randomForest: randomForest( dat.s, vari ~ Depth+slope+aspect+eastness+northness+surf.curv+TPI+TRI+rough+dist.coast+flowdir) | Depth+slope+aspect+eastness+northness+surf.curv+TPI+TRI+rough+dist.coast+flowdir+moran |

921
922    Details on the spatial interpolation methods considered

EVHOE.ARNOGLO.2012

923
924 **Methods using just geographic coordinates or depth (IMGeo)**
925 *1st and 2nd order trend surfaces (LM1 and LM2, respectively).*
926 In these interpolation methods, a first or second order trend surface is fitted to the raw data,
927 respectively. It is a simplistic approach that was included in the current work as a worst case
928 scenario that should be slightly better than a simple overall mean.
929 *Inverse distance weighting (IDW)*
930 Inverse distance weighting (IDW) is an advanced nearest neighbor approach that allows
931 including more observations than only the nearest points. The value at a certain grid cell is
932 obtained from a linear combination of the surrounding locations and the weight of each
933 observation is determined by the distance. IDW is an exact interpolator. The method is fast,
934 easy to implement and easily "tailored" for specific needs, but ancillary data cannot be
935 incorporated. The method tends to generate "bull's eye patterns" (Sluiter, 2009).
936 *Voronoi tessellation (VorT)*
937 The nearest neighbors method predicts the value of an attribute at an unsampled point based on
938 the value of the nearest sample by drawing perpendicular bisectors between sampled points,
939 forming such as Voronoi polygons (or Dirichlet/ Thiessen). This produces one polygon per
940 sample and the sample is located in the center of the polygon, such that in each polygon all
941 points are nearer to its enclosed sample point than to any other sample points (Legendre and
942 Legendre, 1998; Li and Heap, 2008; Webster and Oliver, 2007). This technique is generally
943 used with point data or categorical variables, but can also be used with densities/biomasses
944 (Baddeley et al., 2006; Dauvin et al., 2004; Morfin et al., 2016; Thorson et al., 2015; Zuur et
945 al., 2007).
946 *Thin Plate Splines (TPS)*
947 Thin plate smoothing splines (TPS), formally known as "laplacian smoothing splines". Similar
948 to the previous method, splines are fitted to the sampled data, but in this method, the smoothing
949 parameter is calculated by minimizing the generalized cross validation function (GCV). This

950　method is relatively robust because the minimization of GCV directly addresses the predictive
951　accuracy and is less dependent on the veracity of the underlying statistical model (Hutchinson,
952　1995) (Li and Heap, 2008). We applied this method using package::fields.

953　*Generalized Additive Models (GAM)*
954　 Generalized additive models (GAM) are a semiparametric extension of generalized linear
955　models (GLM), but allow nonlinear relationships between the response and explanatory
956　variables (Wood, 2006), are very commonly used in biological studies (Guisan et al., 2002).
957　GAMs have been often used as a method to produced spatial predictions (i.e. interpolation) by
958　considering the geographic coordinates and its interaction as covariates (Augustin et al., 2013).
959　In the current work we used a GAM applied with the geographic coordinates as covariates
960　(s(x,y, bs="ts"))(GAM), a model where besides geographic coordinates, depth was also
961　considered as covariate (GAMd). GAMs were applied using the r package::mgcv (Wood,
962　2006).

963　*Kriging*
964　From an interpolation point of view, kriging is equivalent to a thin-plate spline and is one
965　species among the many in the genus of weighted inverse distance methods, albeit with
966　attractive properties. However, from a statistical point of view, kriging produces the "best linear
967　unbiased prediction" for an unknown location. It is linear since the estimated values are
968　weighted linear combinations of the available data, unbiased because the mean of the error is
969　0, and it aims to minimize the variance of the errors (Cressie 1990). Several variations of kriging
970　methods were selected following previous works, all applied in the log transformed data. A
971　sequence of interpolation approaches was considered, starting with ordinary kriging with global
972　mean (Okr using automatic modelling and Mkr using manual variogram fitting), ordinary
973　kriging with local neighbourhood estimation (considering 3, 5, 7, 10, 20 and 30 neighbours.
974　Kriging neighborhood is a defined area, in terms of shape and size. Only samples from this area
975　are used in the computation of the local estimates using the kriging technique.
976　Kriging with external trend, also called universal kriging using depth as covariate (Ukr). It is
977　an extension of OK by incorporating the local trend within the neighbourhood search widow as
978　a smoothly varying function of the coordinates. UK estimates the trend components within each
979　search neighbourhood window and then performs SK on the corresponding residuals.

980　*Stochastic conditional simulations (CSim)*
981　These techniques are used more and more, commonly to generate a series of spatial data that
982　have a given degree of spatial dependence, in order to evaluate whether or not observed sample
983　data show significant spatial patterns (Fortin and Dale, 2005). In this method, the parameters
984　of the variogram model (defined previously for ordinary kriging) derived from the experimental
985　variogram were used to generate 1000 stochastic simulations, with the same degree of spatial
986　variance as the observed data (Fortin and Dale, 2005). These methods are known to generate
987　maps having more spatial variability than the kriged ones and hence looking more realistic in
988　comparison to the observed map (Fortin and Dale, 2005).

989　**Methods using topographic covariates (IMCov)**
990　In a mixed method approach to interpolation, the final predictions result of a combination of
991　methods. The main trends are modelled in function of a group of selected covariates in first step
992　(for example General Linear Model (GLM) or machine learning). In a second step, the residuals
993　of this model are then analyzed using kriging, and then incorporated into the predictions (Hengl
994　et al., 2007, 2004; Li et al., 2016). These methods require the availability of covariates. In the
995　current work we used the marine topographic variables derived from GEBCO bathymetric
996　maps, therefore widely available at a worldwide scale. Three mixed models were considered,
997　one using general linear model and two using machine learning algorithms, regression trees and
998　random forest.

999　*Multiple regression (GLM)*

1000 Stepwise multiple regression with all topography-depth covariates was carried out for each
1001 distribution (lvari ~ Depth + slope + rough + moran + TRI + TPI + dist.coast + flowdir + aspect
1002 + eastness + northness). The regression model considered assumes that the residuals are
1003 generated from a normally distributed, second-order stationarity random process—i.e. a random
1004 process that has a constant mean and variance.
1005 *Regression trees (RTre)*
1006 The regression tree approach (also known as binary decision trees) uses binary recursive
1007 partitioning whereby the data of the primary variable are successively split along the gradient
1008 of the explanatory variables into two descendent subsets or nodes. These splits occur so that at
1009 any node the split is selected to maximize the difference between two split groups or branches.
1010 The mean value of the primary variable in each terminal node can then be used to map the
1011 variable across the region of interest (Li and Heap, 2008). Regression tree (CART) algorithm
1012 was fitted to the data to produce a tree with optimal tree size.
1013 *Random forest (RFor)*
1014 A random forest model of each species distribution in function of all marine topographic
1015 covariates was produced (Hengl et al., 2015; Li et al., 2016, 2013, 2011)(vari ~ Depth + slope
1016 + aspect + eastness + northness + surf.curv + TPI + TRI + rough + moran + dist.coast + flowdir).
1017
1018

1019 **Supplement 4**

1020
1021 Correlation plot between the indicators



1022

1023 **Supplement 5**

1024
1025
1026 # title: "Example script for the selection of interpolation method"
1027 # author: "Marta M Rufino^[EMH, IFREMER]"
1028
1029 # The aim of this script is to provide an example of the selection of an interpolation method.
1030 This is an accompanying work of the paper.

```
1031
1032    # We will need to have theese packages intalled:
1033    require(ggplot2); require(RColorBrewer); require(gridExtra) # ploting
1034    require(dplyr); require(tidyr);#data manipulation
1035    require(raster); require(rasterVis); #plot and spatial data manipulation
1036    require(sp); # spatial data
1037    require(gstat); #kriging and idw
1038    require(fields) #tps
1039    require(mgcv); #gam
1040    require(RGeostats) #spatial indicators. Pkg can be download from here
1041    http://rgeostats.free.fr/download.php and install manually
1042    require(ineq) #Gini index (spatial indicators)
1043    require(ggrepel) # PCA plot
1044
1045    # Get an example running with MEUSE dataset:
1046    library(sp); library(gstat)
1047    data(meuse)
1048    data(meuse.grid)
1049    gridded(meuse.grid) = ~x+y
1050    # m <- vgm(.59, "Sph", 874, .04)
1051
1052    dat = meuse %>%
1053      dplyr::select(x,y,zinc, dist) %>%
1054      dplyr::rename("vari" = "zinc",
1055               "Depth" = "dist")
1056    dat.grid <- meuse.grid["dist"]
1057    names(dat.grid) = "Depth"
1058
1059    # 1. Interpolate the data (function interp.dat_CV)
1060
1061    # First we make the different spatial interpolation models. For each model, we will do cross
1062    validation. We will only consider the models without covariates to facilitate the process.
1063
1064    ## This chunk runs a function to interpolate the data of each species distribution
1065    (interp.dat_CV) and estimate respective CV
1066
1067    ## DataDir should the the directory where you have your functions file
1068    ('interp.dat_CV_script.r')
1069
1070    # Please run the required functions which are in the end of the script
1071
1072    # Run the function
1073    kk <- interp.dat_CV(nam="zinc", dat=dat, dat.grid=dat.grid, CV=TRUE, plotit=FALSE,
1074              replicate.cv = 10)
1075    # note we only used 10 replicates of the cv instead of 100 to make a quicker test
1076
1077    # see the results:
1078    head(kk)
1079    # the result of this function is a list with three items:
1080    # 1. res is the raster stack with all predictions from each method;
```

```
1081    # 2. cv.results is the cross-validation summary results
1082    # 3. model.params is some of the models parameters stored
1083
1084    # Extract the raster stack with all interpolator predictions and store it as a new object
1085    pred = kk$res
1086    names(pred)
1087
1088    # Plot the interpolations
1089    coli <- function (region = rev(brewer.pal(n = 10, 'Spectral')), ...)
1090    {theme <- rasterTheme(region = region, ...); theme}
1091    levelplot(pred, par.settings = coli)
1092
1093    # Reshape the table to fit nicely in the results
1094    cv.res <- left_join(
1095      data.frame(kk$cv.results) %>%
1096        dplyr::select(Index, median, method) %>%
1097        tidyr::spread(Index, median),
1098      data.frame(kk$cv.results) %>%
1099        dplyr::filter(Index=="VEcv") %>%
1100        dplyr::select(VEcv.Q1, VEcv.Q3, method) )
1101
1102    # Reorder factor levels
1103    cv.res$method <- factor(cv.res$method,
1104                    levels=c("tre2","GAMp","IDWo",
1105                        "TPSp", "MKri","UKri"))
1106
1107    # Round
1108    cv.res[,-c(1)] <- round(cv.res[,-c(1)],2)
1109
1110    # Make the log of the measures
1111    cv.res$lMAE <- log1p(cv.res$MAE)
1112    cv.res$lRMSE <- log1p(cv.res$RMSE)
1113
1114    # Classify VEcv
1115    cv.res$VEcv.class <- cv.res$VEcv
1116    cv.res$VEcv.class <- cut(cv.res$VEcv.class, c(-2,0,10, 30, 50, 80, 100))
1117    levels(cv.res$VEcv.class) <-
1118      c("0. worst_then_mean",
1119        "1.very_poor",
1120        "2.poor",
1121        "3.average",
1122        "4.good",
1123        "5.excellent")
1124
1125    # Inverted VEcv, i.e. the bigger the worst:
1126    cv.res$VEcv.inv <- abs(cv.res$VEcv/100-1)*100
1127
1128    # Estimate spatial indicators
1129
```

```
1130   # In this chunk we will estimate the spatial indicators using the sampled data and the
1131   interpolated data (prediction rasters).
1132
1133   # For this we will use the pre-packed functions in the packages 'Rgeostats' and 'ineq', although
1134   the indexes are relatively simple to calculate.
1135
1136   # Further, the function will also estimate the 'data limits integrity indicators'.
1137
1138   # This chunk estimates the difference in the spatial indicators between the raw data and
1139   interpolated surfaces and the data limits integrity indicators
1140
1141   # Note we require RGeostats and ineq for this chunk.
1142
1143   # Test the function in one case
1144   fun.inter2(ii="tre2", dat=dat, pred=pred)
1145
1146   # Apply to all interpolation methods:
1147   ind.res <- lapply(as.list(levels(cv.res$method)), fun.inter2, dat=dat, pred=pred)
1148   ind.res <- do.call("bind_rows", ind.res)
1149
1150   # Merge the results with cv results:
1151   tot.res <- full_join(cv.res, ind.res, by=c("method"))
1152
1153   # Now, all the indicators were estimated for each interpolation method. We shall then proceed
1154   to make the first step of the selection method.
1155
1156   # First selection step: VEcv interquantile range
1157
1158   # Reorder factor levels:
1159   tot.res$method <- factor(tot.res$method,
1160                   levels=c("tre2","GAMp","IDWo",
1161                        "TPSp", "MKri","UKri"))
1162   tot.res$short.method <- ordered(tot.res$method,
1163   labels=c("LM2","GAM","IDW","TPS","MKr","UKr"))
1164
1165
1166   # Which methods have the VEcv higher than the lower Q3
1167   tot.res <- tot.res %>%
1168     dplyr::select(-MAE, -R2, -RMSE) %>%
1169     # valid methods, i.e. within range of inter-quartile:
1170     dplyr::mutate(
1171       VEcv.criteria = c(VEcv.Q3 >= max(VEcv.Q1)))
1172
1173   # Plot the VEcv, Q1 and Q3 and respective criteria
1174   tot.res %>%
1175     ggplot(aes(x=method, y=VEcv))+
1176     geom_point(aes(col=VEcv.criteria))+
1177     geom_crossbar(aes(ymin=VEcv.Q1, ymax=VEcv.Q3,y=VEcv, col=VEcv.criteria), alpha=.5,
1178   width=.5)+
1179     ggtitle("First step of interpolators selection")
```

```
1180
1181
1182     # Second selection step: indicators
1183
1184     # select the data for the PCA of indicators
1185     row.names(tot.res) = tot.res$short.method
1186     sec.res <- tot.res %>%
1187       dplyr::filter(VEcv.criteria==TRUE) %>%
1188       dplyr::select(VEcv.inv, lMAE, lRMSE,
1189               CGdist, a.linertia, a.iso, a.Gini,
1190               a.pix_under, a.pix_over, a.mean_perc, a.over_perc)
1191
1192     # Use the function to plot the results and estimate the best method of the selection
1193     PCbiplot(datpc=sec.res, x="PC1", y="PC2")
1194
1195
1196
1197     ########################
1198     ## Functions required
1199     ########################
1200
1201
1202
1203     interp.dat_CV <- function(nam, dat, dat.grid,
1204                     CV=TRUE, plotit=TRUE,
1205                     replicate.cv = 10){
1206       require(ggplot2); require(RColorBrewer);  # ploting
1207       require(dplyr); require(tidyr);#ploting and manover
1208       require(raster); require(rasterVis); #plot and manipulation
1209       require(sp); # spatial data
1210       require(gstat); #kriging and idw
1211       require(fields) #tps
1212       require(mgcv); #gam
1213       #require(scales) #modeling
1214       theme_set(theme_bw(base_size = 9));
1215
1216       # Arguments:
1217       # nam is the label code
1218       # dat data frame with x, y, vari (variable of interest) and Depth
1219       # dat.grid # predictions grid that we want to estimate. class SpatialPixels - sp
1220       # plotit: produce plots for each interpolation.
1221
1222       # Start the function
1223
1224       # make a raster stack to fill with interpolation predictions of the different models:
1225       dat.pred <- raster(dat.grid)
1226       dat.pred[] <- NA
1227       dat.pred <- stack(dat.pred)
1228       # make a dataframe to fit in the parameters
1229       model.params <- data.frame(code=as.character(nam))
```

```r
1230
1231       ## make the spatial object
1232       dat.s <- dat
1233       coordinates(dat.s) <- ~x+y
1234       # the warning is due to the recent change to PROJ6
1235       proj4string(dat.s) <- CRS("+init=epsg:28992")
1236       proj4string(dat.grid) <- CRS("+init=epsg:28992")
1237       #dat.border <- spTransform(dat.border, utm30)
1238       #dat.pred <- stack(raster(dat.grid)) # obj to save the data
1239
1240       # for plotting
1241       coli <- function (region = rev(brewer.pal(n = 10, 'Spectral')), ...)
1242       {theme <- rasterTheme(region = region, ...); theme}
1243
1244       # function to make individual plots
1245       fun.plot <- function(ii, dat.s){
1246        print(levelplot(dat.pred[[ii]]+.1, main=paste(ii, nam, round(max(dat$vari))),
1247                  zscaleLog=FALSE,
1248                  par.settings = coli))}
1249
1250       # Function to estimate error measures
1251       fun.eval <- function(observed, predicted){
1252        resi <- c(observed- predicted)
1253        # rmse(sim=predicted, obs=observed)
1254        (RMSE <- sqrt(mean(resi^2)))
1255        #mae(sim=predicted, obs=observed)
1256        (MAE <- mean(abs(resi)))
1257        #(RMAE = MAE/mean(observed))
1258        #RMAE2 = mean(abs((kk$predicted-kk$observed)/mean(kk$observed)))*100
1259        #(RRMSE = RMSE/mean(observed))
1260        #RRMSE2 = sqrt(mean((kk$predicted-kk$observed)/mean(kk$observed)^2))*100
1261        # R2 should be 1-sum((kk$observed-kk$predicted)^2)/sum((kk$observed-
1262       mean(kk$observed))^2)
1263        (R2 <- 1-(sum((resi)^2)/sum((observed-mean(observed))^2)))
1264        #1-(RMSE/sqrt(mean((kk$observed-mean(kk$observed))^2))))
1265        #(R3 <- 1-var(resi)/var(kk$observed)) #HENGL
1266        (VEcv <- (1 - sum((resi)^2)/
1267             sum((observed-mean(observed))^2))*100)
1268        res.error <- data.frame(RMSE=round(RMSE,2),
1269                   MAE=round(MAE,2),
1270                   #RMAE=round(RMAE,2), RRMSE=round(RRMSE,2),
1271                   R2=round(R2,2),
1272                   VEcv=round(VEcv,2))
1273        return(res.error)
1274       }
1275
1276       # Function to make the cross validation and estimate error measures
1277       cvfun.replicate <- function(xx, FUN, ii=ii, nam=nam, replicate.cv=replicate.cv){
1278        # xx is the data frame with x,y and biom..., FUN is the fun model of each method;
1279        # xx=dat.s; FUN=cv1.fun.cv; replicate.cv=10
```

```
1280        cv2.fun.fold <-function(xx, FUN){
1281          set.seed(seed <- as.integer(runif(1)*2e9))
1282          print(seed)
1283          kf <- sample(rep(seq_len(10), length.out=nrow(dat)))
1284          # Apply fun for the 10 folds
1285          kk <- lapply(as.list(sort(unique(kf))),
1286                       FUN = FUN, xx=xx, kf=kf) %>%
1287            dplyr::bind_rows()
1288          kk$seed=seed
1289          ## if we want to export predicted/observed
1290          #write.table(kk, append=TRUE,
1291          #          file = paste0(paste("pred.obs_cv1000", nam, ii, sep="_"), ".xls"),
1292          #          sep="\t", row.names=FALSE, col.names=FALSE)
1293          assign("last.warning", NULL, envir = baseenv())
1294
1295          kk <- kk %>%
1296            dplyr::group_by(fold) %>%
1297            do(fun.eval(observed=.$observed, predicted=.$predicted)) %>%
1298            dplyr::filter(is.finite(VEcv)) %>%
1299            ungroup() %>%
1300            dplyr::summarise_all(funs(mean)) %>%
1301            dplyr::select(RMSE:VEcv) %>%
1302            data.frame()
1303          return(kk)
1304          rm(kf, kk, seed)
1305        }
1306
1307      # to test cv2.fun.fold(FUN = cv1.fun.cv, xx=dat)
1308      # replicate CV 100 times
1309      xx1 <- replicate(replicate.cv, cv2.fun.fold(FUN = FUN, xx=xx), simplify = FALSE) %>%
1310        bind_rows %>%
1311        mutate(sp=nam, method=ii)
1312      xx1[mapply(is.infinite, xx1)] <- NA
1313      xx1 <- na.exclude(xx1)
1314      ## plot the distribution
1315      #print(xx1 %>% tidyr::gather(Index, value, RMSE:VEcv) %>%
1316      #      ggplot(aes(x=value, group=Index, col=Index))+geom_density()+facet_wrap(~Index,
1317  scales="free"))
1318      ## if we want to export the results:
1319      # write.table(xx1, file = paste("indices_cv1000", nam, ii, ".xls", sep="_"), sep="\t",
1320  row.names=FALSE)
1321
1322      # Get stats
1323      kk1 <- xx1 %>% tidyr::gather(Index, value, RMSE:VEcv) %>%
1324        dplyr::group_by(Index) %>%
1325        dplyr::summarize(VEcv.Q1=quantile (value, probs=0.25),
1326                    VEcv.Q3=quantile(value, probs=0.75),
1327                    mean=mean(value, na.rm=TRUE), N=n(),
1328                    median=median(value, na.rm=TRUE), N=n(),
1329                    max=max(value, na.rm=TRUE),
```

```r
1330                   min=min(value, na.rm=TRUE)) %>%
1331         dplyr::mutate(sp=nam, method=ii)
1332        return(kk1); rm(xx)
1333      }
1334
1335
1336      ####################################
1337      # 2nd order trend surface
1338      ####################################
1339      ii <- "tre2"
1340      dat.trend2 <- gstat(formula=vari ~ 1, data=dat.s, degree=2)
1341      dat.trend2 <- predict(dat.trend2, newdata=dat.grid)
1342      #spplot(dat.trend2[1], contour=TRUE,main="2nd order trend surface interpolation")
1343      dat.pred[[ii]] <- raster(dat.trend2[1])
1344
1345      # Cross validation replicate
1346      if(CV==TRUE){
1347        # function to do CV on each fold
1348        cv1.fun.cv = function(xx, k, kf){
1349          # Function to reproduce the interpolator
1350          # for a part of the data and predict with the other part
1351          # the output MUST be a dataframe with:
1352          # fold/observed/predicted
1353          kk <- gstat(formula=vari ~ 1, data=xx[kf != k,], degree=2)
1354          kk1 <- predict(kk, newdata=xx[kf == k,])
1355          return(data.frame(fold = k, observed = xx[kf == k,]$vari,
1356                     predicted = kk1$var1.pred))
1357          rm(kk, kk1, k)
1358        }
1359        kk <- cvfun.replicate(xx=dat.s, FUN=cv1.fun.cv, ii=ii, nam=nam, replicate.cv=replicate.cv)
1360        cv.results <- kk;
1361      }
1362
1363
1364      # ####################################
1365      # # GAM model
1366      # ####################################
1367      require(mgcv)
1368      ii="GAMp"
1369      dat.mod <- gam(vari~s(x,y, bs="ts"), data=dat)
1370      kk <- data.frame(coordinates(dat.grid));
1371      names(kk)= c("x","y")
1372      dat.mod2 <- predict(dat.mod, newdata=kk)
1373      kk <- cbind(kk, dat.mod2)
1374      dat.pred[[ii]] <- rasterFromXYZ(kk)
1375      # store parameters
1376      model.params$R2.GAMp <- summary(dat.mod)$r.sq
1377
1378      # Cross validation replicate
1379      if(CV==TRUE){
```

```
1380        # function to do CV on each fold
1381        cv1.fun.cv = function(xx, k, kf){
1382          kk <- gam(vari ~ s(x,y, bs="ts", k=50), data=xx[kf != k,])
1383          kk1 <- predict(kk, newdata=xx[kf == k,])
1384          return(data.frame(fold=k, observed=xx[kf == k,]$vari,
1385                    predicted=kk1))
1386          # kk <- gam(vari ~ s(x, y, bs="ts", k=50), data=xx[kf != k,])
1387          # #print(summary(kk)); print(plot(kk))
1388          # kk1 <- predict(kk, newdata=xx[kf == k,])
1389          # return(data.frame(fold=k, observed=xx[kf == k,]$vari, predicted=kk1))
1390          rm(kk, kk1, k)
1391        }
1392        # test: cv1.fun.cv(xx=dat, k=1, kf=kf)
1393        kk <- cvfun.replicate(xx=dat, FUN=cv1.fun.cv, ii=ii, nam=nam, replicate.cv=replicate.cv)
1394        print(head(kk))
1395        cv.results <- bind_rows(cv.results, kk);
1396        rm(kk)
1397      }
1398
1399      rm(dat.mod, dat.mod2, ii)
1400
1401
1402      # ####################################
1403      # # Inverse distance weighting interpolation OPTIMIZED
1404      # ####################################
1405      ii="IDWo"
1406      RMSE <- function(observed, predicted) {
1407        sqrt(mean((predicted - observed)^2, na.rm=TRUE))}
1408
1409      f1 <- function(x, test, train) {
1410        nmx <- x[1]
1411        idp <- x[2]
1412        if (nmx < 1) return(Inf)
1413        if (idp < .001) return(Inf)
1414        m <- gstat(formula=vari~1, locations=train, nmax=nmx, set=list(idp=idp))
1415        p <- predict(m, newdata=test, debug.level=0)$var1.pred
1416        RMSE(test$vari, p)
1417      }
1418      # set.seed(20150518)
1419      i <- sample(nrow(dat.s), 0.2 * nrow(dat.s))
1420      tst <- dat.s[i,]
1421      trn <- dat.s[-i,]
1422      opt <- optim(c(8, .5), f1, test=tst, train=trn)
1423
1424      dat.idwopt <- gstat(formula=vari~1, locations=dat.s, nmax=opt$par[1],
1425    set=list(idp=opt$par[2]))
1426      dat.idwopt <- raster::interpolate(raster(dat.grid), dat.idwopt)
1427      ## [inverse distance weighted interpolation]
1428      dat.idwopt <- mask(dat.idwopt, dat.grid)
1429      dat.pred[[ii]] <- dat.idwopt
```

```
1430
1431       # Cross validation replicate
1432       if(CV==TRUE){
1433         # function to do CV on each fold
1434         cv1.fun.cv = function(xx, k, kf){
1435           kk <- gstat(formula=vari~1, locations=xx[kf != k,], nmax=opt$par[1],
1436    set=list(idp=opt$par[2]))
1437           kk1 <- predict(kk, newdata=xx[kf == k,])
1438           return(data.frame(fold=k, observed=xx[kf == k,]$vari,
1439                       predicted=kk1$var1.pred))
1440           rm(kk, kk1, k)
1441         }
1442         # test: cv1.fun.cv(xx=dat.s, k=1, kf=kf)
1443         kk <- cvfun.replicate(xx=dat.s, FUN=cv1.fun.cv, ii=ii, nam=nam, replicate.cv=replicate.cv)
1444         print(head(kk))
1445         cv.results <- bind_rows(cv.results, kk); rm(kk)
1446       }
1447       rm(i,tst, trn, opt,f1)
1448
1449
1450       ####################################
1451       # TPS
1452       ####################################
1453       print(ii <- "TPSp")
1454       kk <- Tps(coordinates(dat.s), dat.s$vari)
1455       dat.tps <- raster::interpolate(raster(dat.grid), kk)
1456       dat.tps <- mask(dat.tps, dat.grid)
1457       dat.pred[[ii]] <- dat.tps
1458
1459       # Cross validation replicate
1460       if(CV==TRUE){
1461         # function to do CV on each fold
1462         cv1.fun.cv = function(xx, k, kf){
1463           kk <- Tps(coordinates(xx[kf != k,]), xx[kf != k,]$vari)
1464           kk1 <- predict(kk, coordinates(xx[kf == k,]))
1465           return(data.frame(fold=k, observed=xx[kf == k,]$vari,
1466                       predicted=c(kk1)))
1467           rm(kk, kk1, k)
1468         }
1469         # test: cv1.fun.cv(xx=dat.s, k=1, kf=kf)
1470         kk <- cvfun.replicate(xx=dat.s, FUN=cv1.fun.cv, ii=ii, nam=nam, replicate.cv=replicate.cv)
1471         print(head(kk))
1472         cv.results <- bind_rows(cv.results, kk); rm(kk)
1473       }
1474
1475
1476       ####################################
1477       # Manual kriging
1478       ####################################
1479       print(ii <- "MKri")
```

```
1480      dat.vgm <- variogram(vari~1, dat.s)
1481
1482      kk = vgm(psill = max(dat.vgm$gamma), model="Sph", range=max(dat.vgm$dist)/2,
1483    nugget=min(dat.vgm$gamma))
1484      dat.fit <- fit.variogram(dat.vgm, model = kk)
1485      ## plot variogram with respective model
1486      #plot(dat.vgm, dat.fit)
1487
1488      dat.krige <- krige(vari ~ 1, dat.s, dat.grid, model = dat.fit)
1489      #spplot(dat.krige[1])
1490      dat.pred[[ii]] <- raster(dat.krige[1])
1491      model.params <- cbind(model.params, "MKri"=data.frame(nug=dat.fit[1,2], sill=dat.fit[2,2],
1492    range=dat.fit[2,3]))
1493
1494      # Cross validation replicate
1495      if(CV==TRUE){
1496        # function to do CV on each fold
1497        cv1.fun.cv = function(xx, k, kf){
1498          kk <- krige(vari ~ 1, xx[kf != k,], xx[kf == k,], model = dat.fit)
1499          return(data.frame(fold=k, observed=xx[kf == k,]$vari,
1500                       predicted=kk$var1.pred))
1501          rm(kk, kk1, k)
1502        }
1503        # test: cv1.fun.cv(xx=dat.s, k=1, kf=kf)
1504        kk <- cvfun.replicate(xx=dat.s, FUN=cv1.fun.cv, ii=ii, nam=nam, replicate.cv=replicate.cv)
1505        print(head(kk))
1506        cv.results <- bind_rows(cv.results, kk);
1507        rm(kk)
1508      }
1509
1510
1511      ###################################
1512      # Universal kriging
1513      ###################################
1514      print(ii <- "UKri")
1515      dat.vgm <- variogram(vari~Depth, dat.s)
1516
1517      kk = vgm(psill = max(dat.vgm$gamma), model="Sph", range=max(dat.vgm$dist)/2,
1518    nugget=min(dat.vgm$gamma))
1519      dat.fit <- fit.variogram(dat.vgm, model = kk)
1520      #plot(dat.vgm, dat.fit)
1521
1522      dat.krige <- krige(vari ~ 1, dat.s, dat.grid, model = dat.fit)
1523      #spplot(dat.krige[1])
1524      dat.pred[[ii]] <- raster(dat.krige[1])
1525      model.params <- cbind(model.params, "UKri"=data.frame(nug=dat.fit[1,2], sill=dat.fit[2,2],
1526    range=dat.fit[2,3]))
1527
1528      # Cross validation replicate
1529      if(CV==TRUE){
```

```
1530        # function to do CV on each fold
1531        cv1.fun.cv = function(xx, k, kf){
1532          kk <- krige(vari ~ Depth, xx[kf != k,], xx[kf == k,], model = dat.fit)
1533          return(data.frame(fold=k, observed=xx[kf == k,]$vari,
1534                      predicted=kk$var1.pred))
1535          rm(kk, kk1, k)
1536        }
1537        # test: cv1.fun.cv(xx=dat.s, k=1, kf=kf)
1538        kk <- cvfun.replicate(xx=dat.s, FUN=cv1.fun.cv, ii=ii, nam=nam, replicate.cv=replicate.cv)
1539        print(head(kk))
1540        cv.results <- bind_rows(cv.results, kk);
1541        rm(kk)
1542      }
1543
1544
1545      ###################################
1546      ## FINAL PLOTS
1547      ###################################
1548
1549      ## exclude Depth layer and project if wanted
1550      # res <- projectRaster(dat.pred, crs=myCRS)
1551      res = stack(dat.pred[[-1]])
1552
1553      if(plotit == TRUE){
1554        par(ask=TRUE)
1555        #samples
1556        p1<- qplot(data=dat, x=x, y=y, size=vari, col=vari, alpha=.0)+
1557          ggtitle(paste(nam,
1558                  "max:", round(max(dat$vari)),
1559                  "mean:", round(mean(dat$vari))))
1560
1561        # maps
1562        p2 <- levelplot(res,
1563                  main=paste(nam, round(max(dat$vari))),
1564                  zscaleLog=FALSE,layout=c(6, 1),
1565                  par.settings = coli)
1566
1567        # log maps
1568        p3 <- levelplot(res+.1, main="Log", zscaleLog=TRUE,
1569                  layout=c(6, 1),par.settings = coli)
1570        # scaled maps
1571        p4 <- levelplot(scale(res), main=paste('Scaled', nam),
1572                  layout=c(6, 1),par.settings = coli)
1573
1574        # histogram
1575        p5 <- histogram(res,
1576                  xlim=c(0,max(dat$vari)))
1577        # density
1578        p6 <- densityplot(res,
1579                  xlim=c(0,max(dat$vari)))
```

```
1580
1581        # boxplot
1582        p7 <- bwplot(res)
1583
1584        grid.draw(grid.arrange(p2,p3,p4,
1585                        layout_matrix = rbind(c(1,1,1),c(2,2,2),c(3,3,3))))
1586        grid.draw(grid.arrange(p5,p1,p6,p7,
1587                        layout_matrix = rbind(c(1,1,1),c(2,3,4))))
1588        par(ask=FALSE)
1589      }
1590
1591      return(list(res = res,
1592                cv.results = cv.results,
1593                model.params = model.params))
1594
1595      assign("last.warning", NULL, envir = baseenv())
1596      # res <- tidyr::gather(data.frame(res), method, pred, -x,-y)
1597    }
1598
1599    # to test
1600    # kk <- interp.dat_CV(nam="zinc", dat=dat, dat.grid=dat.grid, CV=TRUE, plotit=TRUE)
1601
1602
1603    # Function to estimate the spatial indicators:
1604    fun.indicators <- function(dat) {
1605      ## function to be used in this chunk:
1606      require(RGeostats);require(ineq) # For Gini index
1607      dat <- data.frame(dat)
1608      names(dat) <- c("x","y","pred")
1609
1610      if(max(dat$pred, na.rm=TRUE)!=0){
1611        dat$pred[dat$pred<0] = 0 #to avoid issues with center of gravity
1612        kk <- db.create(x1=dat$x, x2=dat$y, z1=dat$pred)
1613        kk2 <- SI.cgi(kk) # for the centre of gravity, inertia and isotropy-
1614        kk5 <- ineq(dat$pred, type="Gini")# gini index
1615        return(data.frame(t(unlist(kk2)[c(1,3,4,5)]), Gini=kk5))
1616      }}
1617
1618    # To test the function
1619    # fun.indicators(dat = rasterToPoints(pred[[1]]))
1620
1621
1622    # Function to apply the indicators function, estimate the difference between sampled and
1623    interpolated and estimate the data limits intergrity measures:
1624    fun.inter2 = function(ii, dat=dat, pred=pred){
1625
1626      predi = rasterToPoints(pred[[ii]])
1627
1628      # Estimate the indicators of interpolated
1629      res2 <- fun.indicators(predi)
```

```
1630
1631     # Estimate the indicators for raw data
1632     res1 <- fun.indicators(dat[, c("x","y","vari")]) # 0.35 sec
1633
1634     # Absulute difference between sampled and interpolated
1635     res <- abs(res1-res2)
1636     names(res)[c(1,2,5)] = paste0("a.",names(res)[c(1,2,5)])
1637     res$method <- ii
1638
1639     # Diff of center of gravity
1640     res$CGdist <- spDistsN1(as.matrix(res1[,c("center1","center2")]),
1641   as.matrix(res2[,c("center1","center2")]))
1642
1643     # Rescale inertia
1644     res$a.linertia = log1p(res$a.inertia)
1645
1646     # Get number of pixels values over max biom pred
1647     ll <- dim(predi)[1] # number of pixels
1648     mx <- max(dat$vari, na.rm=TRUE)
1649     mm <- mean(dat$vari, na.rm=TRUE)
1650
1651     res$a.pix_under <-
1652       abs(ifelse(is.null(dim(predi[predi[,3]<0,])[1]),0,
1653             round(dim(predi[predi[,3]<0,])[1]/ll*100,2)))
1654
1655     res$a.pix_over <-
1656       abs(ifelse(is.null(dim(predi[predi[,3]> mx,])[1]),0,
1657             round((dim(predi[predi[,3] > mx,])[1]/ll)*100,2)))
1658
1659     res$a.mean_perc <- abs(round(c(mean(predi[,3], na.rm=TRUE)- mm)/mm * 100, 2))
1660
1661     res$a.over_perc <- abs(round((max(predi[,3], na.rm=TRUE)-mx)/mx*100,2))
1662
1663     res <- res[,c("method", "CGdist","a.linertia","a.iso","a.Gini",
1664             "a.pix_under","a.pix_over","a.mean_perc","a.over_perc")]
1665     return(res)
1666   }
1667
1668   # Function to make the PCA and estimate the best method according to the indicators
1669   PCbiplot <- function(datpc=sec.res,
1670             x="PC1", y="PC2") {
1671     require(ggrepel)
1672
1673     # exclude indicators with zero only
1674     datpc = datpc[,colSums(datpc)!=0]
1675
1676     # PCA
1677     PC <- prcomp(datpc, scale=TRUE, center=FALSE)
1678     #biplot(PC)
1679     data <- data.frame(winner2=row.names(PC$x),PC$x)
```

```r
1680     datapc <- data.frame(varnames=rownames(PC$rotation), PC$rotation)
1681
1682     mult <- min((max(data[,"PC2"]) - min(data[,"PC2"])/(max(datapc[,"PC2"])-
1683   min(datapc[,"PC2"]))),(max(data[,"PC1"]) - min(data[,"PC1"])/(max(datapc[,"PC1"])-
1684   min(datapc[,"PC1"]))))
1685     datapc <- transform(datapc, v1 = .7 * mult * (get("PC1")),v2 = .7 * mult * (get("PC2")))
1686     dev <- paste0(c(round((((PC$sdev)^2 / sum(PC$sdev^2) )*100))[1:2],"%")
1687
1688     # Get distance to center of each point:
1689     data$dist <- apply(data[,c("PC1","PC2")], 1,
1690                   function(x) {
1691                     (sqrt((x[1] - 0)^2+(x[2]-0)^2))})
1692
1693     # Reverse weights, as the closer to zero the better:
1694     data$dist2 <- 1/data$dist
1695
1696     # Classification and col of criteria
1697     col.ind <- data.frame(nam=row.names(datapc), class=1)
1698     col.ind[col.ind$nam %in% c("lMAE","lRMSE","VEcv.inv"),2]<-"Error";
1699     col.ind[col.ind$nam %in% c("CGdist","a.linertia","a.iso","a.Gini"),2]<- "Spatial";
1700     col.ind[col.ind$nam %in% c("a.pix_under","a.mean_perc","a.over_perc","a.pix_over"),2]<-
1701   "Integrity"
1702     col.ind$nam=factor(col.ind$nam)
1703     col.ind$col=c("#5E4FA2","#3288BD","#66C2A5")[factor(col.ind$class)]
1704     #rev(brewer.pal(11, "Spectral"))
1705
1706     plot1 <-
1707       ggplot(data, aes(x=PC1, y=PC2)) +
1708       geom_point(aes(col=dist2), size = 1, shape=16)+
1709       geom_text_repel(aes(label = winner2, size=dist2, color=dist2)) +
1710       scale_colour_gradient(high = "#9E0142", low = "#FDAE61")+
1711       geom_hline(aes(yintercept=0), size=.2, color=8, linetype=2) +
1712       geom_vline(aes(xintercept=0), size=.2, color=8, linetype=2)+
1713       xlim(extendrange(c(data$PC1,datapc$PC1))[1],0.01)+
1714       ylim((extendrange(c(data$PC2,datapc$PC2))))+
1715       # plot criteria:
1716       geom_text_repel(data=datapc, aes(x=v1, y=v2, label=varnames),size = 3,
1717   segment.alpha=.5,
1718                     color=col.ind$col)+
1719       geom_segment(data=datapc, aes(x=0, y=0, xend=v1, yend=v2),
1720   arrow=arrow(length=unit(0.2,"cm")), color=col.ind$col)+
1721       xlab(paste0("PC1 (",dev[1],")"))+
1722       ylab(paste0("PC2 (",dev[2],")"))+
1723       ggtitle("PCA of indicators")+
1724       theme(line = element_blank(),
1725           axis.text=element_blank(),
1726           axis.ticks=element_blank())+
1727       scale_size(range = c(3, 5))+
1728       guides(size=FALSE, fill=FALSE, col=FALSE)
1729     plot2 <-
```

```r
1730      ggplot(data, aes(x=reorder(winner2,dist2), y=dist2)) +
1731      geom_bar(stat="identity", aes(fill=dist2),col="White", alpha=.8)+
1732      scale_fill_gradient(high = "#9E0142", low = "#FDAE61")+
1733      theme(line = element_blank(),
1734          axis.text.y=element_blank(),
1735          axis.ticks.y=element_blank(),
1736          axis.text.x=
1737            element_text(angle=90,hjust=1))+
1738      xlab("Interpolation methods")+
1739      ylab("Inv. dist. to center")+
1740      ggtitle(" ")+
1741      guides(size=FALSE, fill=FALSE, col=FALSE)
1742
1743    grid.draw(arrangeGrob(plot1, plot2, ncol=2, widths = c(3/4,1/4)))
1744
1745    # Result's table
1746    kk <- data %>%
1747      dplyr::rename("method" ="winner2") %>%
1748      dplyr::mutate(dist2 = round(dist2,2)) %>%
1749      dplyr::select(method, dist2) %>%
1750      dplyr::arrange(dist2)
1751    return(kk)
1752  }
1753


1754

1755
```