

This is the authors' final version of the contribution published as:

Gobeyn, S., Mouton, A.M., **Cord, A.F., Kaim, A., Volk, M.**, Goethals, P.L.M. (2019):
Evolutionary algorithms for species distribution modelling: A review in the context of
machine learning
Ecol. Model. **392**, 179 - 195

The publisher's version is available at:

<http://dx.doi.org/10.1016/j.ecolmodel.2018.11.013>

1 Evolutionary algorithms for species distribution modelling: a review
2 in the context of machine learning

3 Sacha Gobeyn^{a,*}, Ans M. Mouton^a, Anna F. Cord^b, Andrea Kaim^b, Martin Volk^b, Peter L.M.
4 Goethals^a

5 ^a*Ghent University, Department of Animal Sciences and Aquatic Ecology, Coupure Links 653, B-9000 Ghent,*
6 *Belgium*

7 ^b*UFZ - Helmholtz Centre for Environmental Research, Department of Computational Landscape Ecology,*
8 *Permoserstr. 15, 04318 Leipzig, Germany*

9 **Abstract**

10 Scientists and decision-makers need tools that can assess which specific pressures lead to ecosystem
11 deterioration, and which measures could reduce these pressures and/or limit their effects. In this
12 context, species distribution models are tools that can be used to help asses these pressures. Evo-
13 lutionary algorithms represent a collection of promising techniques, inspired by concepts observed
14 in natural evolution, to support the development of species distribution models. They are suited
15 to solve non-trivial tasks, such as the calibration of parameter-rich models, the reduction of model
16 complexity by feature selection and/or the optimization of hyperparameters of other machine learn-
17 ing algorithms. Although widely used in other scientific domains, the full potential of evolutionary
18 algorithms has yet to be explored for applied ecological research. In this synthesis, we study the role
19 of evolutionary algorithms as a machine learning technique to develop the next generation of species
20 distribution models. To do so, we review available methods for species distribution modelling and
21 synthesize literature using evolutionary algorithms. In addition, we discuss specific advantages and
22 weaknesses of evolutionary algorithms and present a guideline for their application. We find that
23 evolutionary algorithms are increasingly used to solve specific and challenging problems. Their
24 flexibility, adaptability and transferability in addition to their capacity to find adequate solutions
25 to complex, non-linear problems are considered as main strengths, especially for species distribution
26 models with a large degree of complexity. The need for programming and modelling skills can be
27 considered as a drawback for novice modellers. In addition, setting values for hyperparameters

*Corresponding author. Tel.: +32 9 264 39 96; Fax: +32 9 244 44 10
Email address: Sacha.Gobeyn@UGent.be (Sacha Gobeyn)

28 is a challenge. Future ecological research should focus on exploring the potential of evolutionary
29 algorithms that combine multiple tasks in one learning cycle. In addition, studies should focus on
30 the use of novel machine learning schemes (*e.g.* automated hyperparameter optimization) to apply
31 evolutionary algorithms, preferably in the context of open science. This way, ecologists and model
32 developers can achieve an adaptable and flexible framework for developing tools useful for decision
33 management.

34 **Keywords**

35 Evolutionary algorithms; machine learning; species distribution modelling; model development; automated
36 model selection

37 **1. Introduction**

38 Innovations in remote sensing and micro-control units used for (near-)real-time monitoring of ecosystems
39 are challenging ecologists to deal with a great number of data coming from different sources (Hampton
40 et al., 2013). Machine learning plays an increasingly important role to deal with this challenge, not only
41 in the field of ecology (Meyers et al., 2017). However, ecologists often struggle with the interpretation of
42 models developed with novel machine learning algorithms (Araújo and Guisan, 2006). As a consequence,
43 scientists are required to search for new approaches to increase reliability, transparency and flexibility of
44 the models they are developing (Elith and Leathwick, 2009).

45
46 The models developed for ecological research and management with machine learning are mainly classified
47 under species distribution models (SDMs). These models aim to define the species-environment relation and
48 from this, estimate the species' geographical distribution. SDMs rely on the concept of an 'ecological niche',
49 described by Hutchinson (1957). This theory conceptualizes the relation between a species' environment
50 and its occurrence (Hirzel and Le Lay, 2008). Niche theory (Hutchinson, 1957) states that a species can only
51 exist if the local combination of environmental gradients, the niche, allows a positive population growth
52 rate, in the absence of immigration. In addition, the theory states that a difference in species' traits allows
53 them to occupy a different niche and coexist in a given spatial unit. To inspect these species' traits, machine
54 learning algorithms are used to classify a species as present, or absent, given the environmental conditions.

55 Often a probabilistic framework is used to express the chance for a species to occur. The environmental
56 conditions are quantified in a number of environmental features, for instance, temperature, precipitation,
57 soil moisture (He et al., 2015). The number of input features, after preprocessing, reach up to 20, while
58 studies using over 30 features are found in literature (Bennetsen et al., 2016). Machine learning is used
59 to train a model estimating a response variable, species occurrence, based on these environmental predictors.

60

61 Species distribution model development with machine learning embeds little ecological hypotheses in the
62 training process as the machine learning algorithms primarily aim to uncover patterns in data (Saeys et al.,
63 2007; Mount et al., 2016). At first instance, this development could be considered as “black box” modelling,
64 which is the case in a number of techniques, *e.g.* artificial neural networks. Yet, there are techniques, such
65 as decision trees, that present interpretable models. In these models, the user can interpret why estima-
66 tions were done in that way. For example, in decision trees, a set of hierarchical rules can be analysed
67 that lead to an estimated species presence. In specific applications, for instance in freshwater management,
68 embedding hypotheses and assumptions is of major importance to preserve ecological interpretability of the
69 developed models (Adriaenssens et al., 2004). Also, incorporating species dispersal and interactions in the
70 next generation of species distribution models (SDMs) requires a hypothesis-driven approach. Evolutionary
71 algorithms (EAs), a collective of machine learning techniques inspired on the concept of evolution, allow
72 embedding these hypotheses by separating model performance evaluation from solution searching (Rauch
73 and Harremoës, 1999). In these EAs, models with parameters and input variables are encoded in so-called
74 ‘chromosomes’. Specific algorithm functions, called genetic operators, are applied to these chromosomes to
75 search for well-performing models.

76

77 Because EAs algorithms offer this transparency and flexibility, we evaluate in this synthesis the role of
78 EAs as a machine learning method for species distribution modelling. We aim to answer the following
79 questions: What is the current role of EAs in species distribution modelling (section 2)?; How do EAs work
80 and what are their specific strengths (section 3)?; What opportunities do these EAs have in the field of
81 species distribution modelling (section 4)?; What guidelines can be given to applying an EA or another
82 metaheuristic as a machine learning method to identify transparent and accurate models (section 5)? It is
83 important to note that this paper focuses on EAs, having in mind that other metaheuristics such as particle
84 swarm optimization, ant colony optimization and simulated annealing also exist. For an overview of these

85 other methods, we refer to supportive information 1. Differences between EAs and other metaheuristics
86 are discussed throughout this manuscript. For an extensive description of EAs and other metaheuristics,
87 we refer readers to Gendreau and Potvin (2010) and Kacprzyk and Pedrycz (2015).

88 **2. Machine learning in species distribution modelling**

89 Many ecological researchers rely on machine learning for the development of SDMs. The use of machine
90 learning has introduced new concepts important for ecologists to understand. In contrast to previous re-
91 views discussing the development of SDMs (*e.g.* Guisan and Thuiller (2005); Araújo and Guisan (2006);
92 Austin (2007)), this review focusses on SDM learning and its technical challenges, rather than development
93 through ecological reasoning. That is why we focus on the discussion of EAs in the context of machine
94 learning.

95

96 Machine learning can be defined as the field of research using computer programs or algorithms that have
97 the ability to adapt or change from an experience (data) given a(n) performance or objective measure.
98 Algorithms to facilitate machine learning are widely applied in many scientific fields such as artificial
99 intelligence, telecommunication and engineering of electronics (web of science, accessed on 13/09/2018).
100 Machine learning algorithms can be used to train models with data so these models can make as good
101 as possible predictions on new, ‘unseen’, data. Machine learning algorithms can be categorized based on
102 whether output labels are (not) used for training, *i.e.* (un-)supervised learning (Box 1). To guide the
103 readers, Box 1 shows a number of definitions used in the field of machine learning, also used in this review
104 paper.

105 In species distribution modelling, supervised learning is typically applied to classify species occurrence in
106 geographical, and possibly the temporal, dimensions. To train SDMs, binary labelled data (species presence
107 or absence) together with environmental input data are used by the machine learning algorithm. In Table 1,
108 an overview of methods used in species distribution modelling are shown, together with a short explanation,
109 the cumulative number of papers mentioning the method (web of science, accessed 08/11/2017), and key
110 references. In addition, in Figure 1, the first report of the method in scientific literature is shown. The
111 remainder of this section aims to shortly introduce these methods to guide readers to the most used ones.
112 Acronyms of these methods can also be found in Table 1.

113

114 Generalized linear and adaptive models (respectively GLMs, and GAMs) and decision trees are the first
115 machine learning techniques used in species distribution modelling (Figure 1). GLMs (and GAMs) are a
116 collection of (semi-)parametric techniques based on three elements: a random component that assumes a
117 probability distribution of a response variable (1), a systematic component specifying linear combination of
118 the explanatory variables with their respective slopes (2) and a link function describing the relation between
119 the random and systematic component (3) (Nelder and Wedderburn, 1972). Decision trees are classifiers
120 expressed as recursive partitions or trees of the feature space (Rokach and Maimon, 2015). These are
121 tree-like representations of a rule-induction, *i.e.* a set of ‘if-then’ rules that are followed leading to a (prob-
122 ability of) species presence or absence. Different algorithms are available to develop decision trees, such
123 as CART (Breiman et al., 1984) and C4.5 (Quinlan Ross, 1993), using a Gini index and entropy measure,
124 respectively. A robust approach based on decision trees is a random forest (RF) (Breiman, 2001). RF uses
125 bootstrap aggregation to generate on a number of decision or regression trees. In bootstrap aggregation,
126 several bootstrap samples from a training data set (objects in the instance space) are taken to develop a
127 number of models. RF is used often in species distribution modelling, and has shown to be an interesting
128 technique to model complex systems, including species interactions (Veza et al., 2015).

129

130 The genetic algorithm for rule set production (GARP) is an EA-inspired method used to produce a rule-
131 bank SDM (Stockwell and Noble, 1992). A rule-bank SDM is a model based on a set of hierarchical rules
132 estimating species presence or absence. In this sense, the model structures obtained with GARP are similar
133 to those obtained with decision trees. GARP is the first SDM software package using EAs, knowing a
134 number of successful applications, for instance, to estimate the effect of global change on species distribu-
135 tions (Peterson et al., 2002). Another technique often used in combination with EAs are artificial neural
136 networks (Ding et al., 2013). Artificial neural networks which are non-linear mapping structures inspired
137 by the biological system of the brain have been successfully used for freshwater applications (Goethals
138 et al., 2007; Muñoz-Mas et al., 2017). From these two examples, it is clear that EAs are often used to
139 support the development of the model structure, *i.e.* search for the most optimal model structure. Despite
140 their success, the popularity of artificial neural networks and also GARP is recently declining compared to
141 that of maximum entropy modelling (Maxent) (Table 1). Introduced in 2006, Maxent uses the principle of
142 maximum entropy to make predictions from incomplete knowledge (Phillips et al., 2006). The maximum
143 entropy method (Maxent) approach uses the principle of maximum entropy to make predictions from in-

144 complete knowledge. The principle of maximum entropy states that the best approximation to an unknown
145 distribution, given a number of constraints, is the distribution which only satisfies these constraints and no
146 others. In other words: Maxent aims to model everything that is known (constraints) but assumes nothing
147 about what is unknown. Maxent is currently the most used package to train SDMs (200+ papers in 2016,
148 based on abstract search, web of science, 08/11/2017). Its theoretical basis, the default use of regularisa-
149 tion (*i.e.* penalize model complexity), flexibility and performance are the main factors explaining Maxent's
150 popularity (Elith et al., 2011). Model complexity can also be penalized in many other algorithms, however,
151 Maxent was the first approach to include it by default. This default inclusion stressed the importance of
152 regularization among users. Peterson et al. (2007) compared the transferability (test on unseen data) of
153 Maxent and GARP and showed that both have their specific advantages. This is not surprising as the
154 'No Free Lunch Theorem' (Wolpert and Macready, 1997) depicts that no algorithm will work well on all
155 problems. A good approach to deal with the 'No Free Lunch Theorem' is to combine and/or compare
156 several approaches in one modelling effort. Thuiller (2003) developed a platform for implementing differ-
157 ent techniques. His initial aim was to present a framework able to simultaneously fit different models to
158 data. It was only six years later that Thuiller et al. (2009) presented a new version of BIOMOD, including
159 the concepts of uncertainty estimation, and ensemble forecasting (Araújo and New, 2007). A weakness
160 of the BIOMOD is that the platform is bound by specific implementations of techniques. Golding et al.
161 (2017) deals with this issue by implementing a modular framework for species distribution modelling. He
162 argues that algorithm success is partly depicted by method transparency empowered by clear encoding and
163 guidelines of use, *i.e.* which method and specific settings are suitable to solve the problem at hand? It is
164 important to note that besides algorithm guidelines, also appropriate data cleansing techniques can consid-
165 erably improve results. An excellent guide for data cleansing of ecological data is given by Zuur et al. (2010).

166

167 Other techniques are available but are not categorized under machine learning as their origins are rooted
168 in niche theory of Hutchinson (1957). These methods, specifically GROWEST (Nix et al., 1977), CLIMEX
169 (Sutherst and Maywald, 1985) and BIOCLIM (Booth, 1985; Nix, 1986) were used in the early days of
170 mapping a species' niche (Figure 1). They are currently less used because of their simplicity, lack of ac-
171 curacy and inability to account for variable interaction (Booth et al., 2014) (Table 1). Another less-used
172 approach is the development of models with fuzzy logic. Fuzzy logic models allow integrating expert knowl-
173 edge in their model structure. Specifically, fuzzy models allow reflecting uncertainty present in linguistic

174 information (Adriaenssens et al., 2004). Although fuzzy logic model development might not be classified
175 under machine learning, it is often used in conjunction with a machine learning algorithm (Chen et al., 2003).

176

177 The methods discussed above generally only consider the species-environment relationship to estimate
178 species occurrence. The spatial structure of the relations is implicitly included in the scale. Spatially
179 explicit models, which incorporate the spatial space in their structure, allow describing processes such as
180 migration and dispersal in a spatial context (DeAngelis and Yurek, 2017; Dunning et al., 1995). Up until
181 today, these explicit methods are less popular, mainly because of their complexity, and need for detailed
182 information to parametrize them (DeAngelis and Yurek, 2017). Given these disadvantages, spatially ex-
183 plicit models do hold a lot of potential to help uncover species behaviour and distribution as a function of
184 environmental pressures.

185

186 Although the role of GARP is recently declining, a number of specific applications of EAs are observed in the
187 literature since 2000 (see Table 2). These novel techniques are mainly applied in the context of freshwater
188 management and used to estimate a link between river modification and freshwater species occurrence.
189 They facilitate feature selection (feature selection, for definition, see Box 1) for artificial neural network
190 and decision tree models (D’heygere et al., 2006) or to estimate model parameters of fuzzy logic SDMs
191 (Fukuda et al., 2011). The results of our literature review (for methodology, see supportive information 2)
192 show that these algorithms are often ‘tailor-made’, and characterized by a specific algorithm formulation.
193 In the next section, we discuss this ‘EA-literature’ with the aim to identify which specific applications and
194 workflows are mainly adopted. Before we do so, we introduce the basic principles of evolutionary algorithms.

195 **3. Evolutionary algorithms**

196 *3.1. Introduction to evolutionary algorithms*

197 EAs aim to solve complex problems by incorporating elements of structured randomness in their search
198 behaviour motivated by principles in evolution such as selection, mutation and crossover (Maier et al.,
199 2014). EAs distinguish from single-point based methods by iterating a population of candidate solutions
200 to an optimum. These solutions are quantified in a fitness mimicking the evolutionary concept.

201

202 EAs have been successfully applied to solve specific problems in water resources management (Maier et al.,
203 2014), astrophysics, bio-informatics (Pal et al., 2006; Sirbu et al., 2010) and software engineering (Eiben
204 and Smith, 2015). Even more, the use of EAs in artificial intelligence is pushing advances in evolving
205 digital objects (software) towards physical embodied artificial evolution (*i.e.* hardware, robots, 3D-printers).
206 Numerous examples exist and the number of applications is expected to increase in the coming years (Eiben
207 and Smith, 2015). EAs have mainly been used as a machine learning method, also to train other methods
208 such as artificial neural networks and decision trees. Their success can be explained by a number of reasons
209 (Eiben and Smith, 2015; Maier et al., 2014):

- 210 1. EAs are assumption-free which make them generally applicable and easily transferable to other prob-
211 lems.
- 212 2. They are flexible and can easily be used in combination with other methods, for instance, other search
213 methods or fuzzy logic.
- 214 3. They are capable to solve complex problems without the need for model simplification often re-
215 quired by traditional optimization methods. Moreover, they are able to uncover less obvious or even
216 unexpected patterns.
- 217 4. The found solutions with EAs allow for an in-depth analysis since a number of near-optimal solutions
218 are generated.

219 EAs iterate a population of chromosomes over a number of generations with genetic operators, *i.e.* selection,
220 crossover and mutation (Box 2 and Figure 2, panel A). This process is inspired by the concept of evolution
221 where genetic information and characteristics in a population are passed on generation by generation. The
222 chromosome is the algorithms' building block storing the formulation and performance of a candidate so-
223 lution to a problem, *i.e.* the genome and fitness. The fitness can be defined as a quantification measure of
224 how good a solution to a problem is. The formulation of the candidate solution is stored in a specific data
225 type, also called genome. For genetic algorithms (GAs), binary or real-valued strings are coded as data type
226 whereas tree-like structures are used for genetic programming (GP) (Figure 2, panel B). Both GAs and GP
227 are classified under evolutionary algorithms. Other types of EAs, such as evolutionary programming and
228 evolution strategies exist (Weise, 2009). GAs differ from other evolutionary algorithms in the way they are
229 designed as problem-independent solvers, whereas other EAs are designed and implemented to solve specific
230 problems. We consider GP to be developed for specific problems. The conceptual difference between GAs
231 and evolutionary programming is that the basic object is considered to be a species rather than a chro-

232 mosome. In evolutionary programming, recombination (crossover) is often not considered. Evolutionary
233 strategies show resemblance to real-valued GAs, but with a focus on the selection and mutation operator.
234 These problem-specific techniques are used in species distribution modelling, but not often. In the next
235 section, we explore the use of this problem-(in)dependent EAs.

236

237 The initialisation of a population with a number of chromosomes (population size, PS) and their genomes
238 is the first step. For binary string GAs, this consists of creating a random string of bits (example lower left
239 panel, Figure 2) with every bit either having the value zero or one. For real-valued strings, a uniform value
240 within a defined interval is chosen for every bit. After initialisation, the fitness is evaluated by mapping
241 the genome to a model with a mapper function. As an example, for feature selection in SDMs, a ‘011’
242 string is translated to the exclusion of the first feature and inclusion of the last two in the model (D’heygere
243 et al., 2006). For parameter estimation, a binary string is translated to an integer or decimal (respectively
244 ‘011’ $\rightarrow (0*1+1*2+1*4) \rightarrow 6$ or $1/6$) for the value of the model parameters (Van Broekhoven et al., 2007).
245 The models are then evaluated with a user-defined objective function and training data leading to a fitness
246 value. Usually, a measure of agreement between the model output and training data is calculated. After
247 fitness evaluation, selection, crossover and mutation operators are applied to the population. The selec-
248 tion operator selects a number of chromosomes from the population as parents to generate offspring based
249 on their fitness value and a selection procedure (*e.g.* tournament selection, roulette wheel selection). In
250 tournament selection, tournaments are organised in which two candidate parents are (randomly) selected
251 and the candidate with the highest fitness is selected as a parent. For roulette wheel selection, parents are
252 selected with a chance proportional to their relative fitness. The crossover operator generates offspring by
253 inheriting a part of the parents’ genomes. For instance, in a GA one-point crossover operator, a position
254 in the genome is randomly chosen as a breakpoint. The parents’ substrings are then combined to form
255 a new string for the offspring (Figure 2, lower panel). A crossover rate determines the probability that
256 crossover between parents occurs. The mutation operator changes the values in random positions in the
257 genomes (or alleles) of the offspring with a rate equal to the mutation rate. After the application of the
258 three operators, the fitness of the new chromosomes is evaluated. Next, a new generation is produced by
259 applying the before-introduced operators. This procedure is repeated until a certain stopping criterion is
260 met. Typically, this criterion is a maximum number of generations or a fitness convergence criterion.

261

262 Besides the context of use of GAs and GP, the way of problem encoding and consequently the implementa-
263 tion of the crossover and mutation operator is different (Mcdermott and O'Reilly, 2015; Rowe, 2015). For
264 GAs, an example of a crossover of two binary strings with a one-point uniform operator is shown in Figure 2
265 (lower left panel). Here, a random number between two (one) and the length (length minus one) of the two
266 parents' genome is chosen and before (after) this position a breakpoint is appointed. The genome for the first
267 offspring is formed by merging the part before and after the breakpoint of parent one and two, respectively.
268 Similarly, for the second offspring the parts before and after the breakpoint are used but now the genomes
269 of parent two and one are used. For GP (Figure 2, lower right panel), breakpoints are chosen between nodes
270 of the tree and these are switched between the parents' genomes. For mutation in a binary string, a random
271 position is chosen and the value for the allele at that position is flipped to the other value (0→1 or 1→0).
272 In case of real-valued strings, a new random value bounded by a predefined interval is chosen at a random
273 position. For tree-like structures, a random terminal or non-terminal node is chosen and replaced with a
274 terminal node or random initiated subtree (see Figure 2, lower right panel) (Mcdermott and O'Reilly, 2015).
275

276 3.2. Application to species distribution modelling

277 In order to obtain an insight into the use of EAs in species distribution modelling, literature abstracts
278 were scanned in the web of science catalogue. The followed methodology to conduct this literature review
279 can be found in supportive information 2. Here, implementations that differ from GARP are discussed,
280 as these implementations vary as a function of the context of the problem. The results of this literature
281 review are shown in Table 2. In this table, we make a clear distinction between 'parameter estimation'
282 and 'hyperparameter optimization'. Parameter estimation refers to the estimation of a unique set of model
283 parameter values (Box 1). With respect to species distribution and ecological modelling, this implies that
284 model parameters that describe the limits of a species' environmental range are estimated, *e.g.* what are
285 the threshold river temperatures in which a fish can survive? Or what are temperature tipping points at
286 which species reproduction declines? As such, parameters are an element of the SDM. Hyperparameter op-
287 timization refers to the search for values of algorithm settings which influence an algorithm's performance.
288 In other words, a hyperparameter can be considered as an algorithms' free option available for the user.
289 The number of neurons and hidden layers are examples of two hyperparameters that need to be set in order
290 to develop an artificial neural network. Or for a RF, on has to set the maximum depth of a tree and the
291 number of trees. As such, values for hyperparameters can be considered as choice elements of the algorithm,

292 and can thus not be directly estimated with data. In the context of Table 2, ‘hyperparameter optimization’
293 in the column ‘subject of training’ refers to the action of using an EA to perform hyperparameter optimiza-
294 tion of another machine learning approach. It is important to differentiate between the column ‘subject of
295 training’ from the columns ‘hyperparameters’ and ‘hyperparameter optimization’, as the latter two refer to
296 (the setting of) hyperparameters of the EA itself.

297

298 Table 2 shows the general characteristics of studies in which EAs are applied: generally they are applied
299 in freshwater management to estimate model parameters, perform feature selection or hyperparameter op-
300 timization of other machine learning techniques. For 14 of the 27, an EA is used for only feature selection
301 whereas in seven studies for only parameter estimation. In five studies, an EA is used for parameter estima-
302 tion and feature selection and in the remaining study, an EA is used for feature selection and optimization
303 of hyperparameters of a decision tree. In case of feature selection, the EAs are used as wrapper methods for
304 other methods; artificial neural networks, *e.g.* D’heygere et al. (2006), and decision trees, *e.g.* Boets et al.
305 (2013). In this approach, the genomes are translated to features for another machine learning technique
306 fitting the response patterns to environmental conditions. In the case of parameter estimation, the EAs are
307 used to estimate the model parameter values of fuzzy logic models, *e.g.* Fukuda (2009); Van Broekhoven
308 et al. (2007). In other words, the model parameters describing suitability range of environmental conditions
309 for a species are estimated.

310

311 A noteworthy observation is that 22 papers presented in Table 2 situate within the domain of freshwater
312 science. The data used in these case studies are often characterized by a high degree of uncertainty and
313 noise, and observation bias, *i.e.* more (less) presence instances are available than absence (see column prev.
314 in Table 2). The latter, causing a bias in model training (Mouton et al., 2010), is the reason why a number
315 of classification measures are typically used in these studies. Often-used measures are listed in Table 3,
316 together with their acronyms. In these measures, species occurrence estimated by the classifier is tested
317 to observations. According to the study objectives, and available data, a set of measures is selected and
318 analysed, each weighting a degree of correct estimation of species presence, on the one hand, and absence,
319 on the other (Mouton et al., 2010). Non-binary measures, such as the (root) mean of squared error, cor-
320 relation and sum of squared errors, are used for regression. In these cases, the probability of occurrence
321 is not estimated, but species numbers (D’Angelo et al., 1995) or density (Fukuda, 2009). In addition, the

322 (root) mean square error and linear correlation is used. In one case, the mean squared error between the
323 non-classified preference (between 0 and 1) and observed presence/absence is computed (Fukuda et al.,
324 2012). In another case, models are penalized for their complexity (Muñoz-Mas et al., 2016a). The trade-off
325 between omission and commission errors are never explicitly considered in model training, although they
326 are considered implicitly by weighting objectives. Assumed prevalence-independent measures like Cohen’s
327 Kappa or the true skill statistic are used to cope with this trade-off, however, there is no agreement whether
328 these are truly prevalence independent (Mouton et al., 2011). In three of the studies reported in Table 2,
329 the training data are stratified by sampling an equal number of presence and absence instances in order to
330 deal with this prevalence dependency (Mouton et al., 2009).

331

332 A number of different implementations of EAs have been used in species distribution modelling. Simple
333 genetic algorithms are generally used and are considered problem-independent. These algorithms apply a
334 GA with uniform crossover and mutation operators, in conjunction with a tournament selection operator
335 (*e.g.* Boets et al. (2013)). Derivative methods have been used in combination with GAs allowing to improve
336 the local search performance of the GAs (Muñoz-Mas et al., 2016b, 2017). In addition, GPs are used, but
337 only limited (Jeong et al., 2011; McKay, 2001; Whigham, 2000). Another interesting application is the
338 use of Bayesian theory in GAs (McClean et al., 2005; Termansen et al., 2006). Feature selection is always
339 implemented in binary strings whereas binary and continuous strings are used for parameter estimation.
340 Crossover rates vary from 0.6 to 0.95 whereas mutation rate are generally lower, between 0.1 and 0.3, with
341 the exception of 0.6 (Muñoz-Mas et al., 2017) and 0.75 (Muñoz-Mas et al., 2016b). A number of 20 to 200
342 chromosomes are reported to iterate over generally 20 to 100 generations. However, a larger number of
343 generations (≥ 1000) are observed in four studies. Selection rates are never reported, as these are typically
344 equal to 50 %.

345

346 Model robustness is tested by applying cross-validation and repeated learning with the same or different
347 samples of the data. Cross-validation is generally used (14 of 27 cases) to test robustness. In this approach,
348 the data are partitioned in a number of samples, i.e. folds. Next, the model is identified with n-1 folds and
349 validated with the remaining fold. In a number of publications, the EA analysis is repeated a number of
350 times with the same data starting from different initial populations in order to test the robustness of the
351 EA (see D’Angelo et al. (1995); Fukuda (2009)). This is because the obtained near-optimal solution might

352 not be equal in every EA run since the search behaviour is characteristic by random choices. An interesting
353 application of this repeated EA analysis is the multilayer perceptron ensembles (a type of artificial neural
354 network) for the modelling of the redfin barbel (Muñoz-Mas et al., 2017). Here, a derivative GA analysis
355 is repeated a number of times to increase ensemble size. By checking convergence of the solutions deter-
356 mined with the GA for an increased ensemble size, one can determine an optimal set of solutions. With
357 this, the authors showed the potential of using multilayer perceptron ensembles and EAs for the identifica-
358 tion of multiple models reflecting simulation uncertainty (i.e. ensemble forecasting (Araújo and New, 2007)).

359

360 An interesting observation is that EA hyperparameters are reported more consistently in recent years. In
361 addition, testing the training robustness as a practice has increased. The latter is probably due to the
362 availability of growing computational resources. In addition, species prevalence is increasingly reported,
363 suggesting that practioners are more aware of the effect of prevalence bias on model training. As such, it is
364 observed that more detailed and robust approaches are presented. Finally, most studies rely on an extension
365 of simple genetic algorithms, whereas genetic programming has not been employed in recent years.

366 **4. Strengths, weaknesses, opportunities and threats analysis**

367 To explore the potential of EAs in species distribution modelling, we performed a strength, weaknesses,
368 opportunities and threats (SWOT) analysis. In this review, strengths and weaknesses refer to current
369 characteristic of EAs that offer respectively advantages or disadvantages compared to other techniques.
370 Opportunities and threats refer to future (dis-)advantages. To perform this analysis, the literature of EAs
371 in species distribution modelling was scanned and the specific strengths and weaknesses of the use of EAs
372 were compiled. In addition, opportunities and threats were assessed by testing compliance with known
373 challenges in species distribution modelling (based on Araújo and Guisan (2006); Austin (2007); Araújo
374 and New (2007); Guisan and Zimmermann (2000); Guisan and Rahbek (2011)). In the past, specific ad-
375 vantages of EAs and metaheuristics were mainly derived from experiments as the true functioning of the
376 algorithms was poorly understood (Boussaïd et al., 2013; Maier et al., 2014). That is why the analysis in
377 this section is based on specific examples rather than theoretical studies.

378

379 EAs are particularly useful in situations where solutions to complex problems have to be found for which
380 little information is available to characterize the optimal solutions. In these situations, it is not possible

381 to do a grid-search, *i.e.* check all candidate solutions one by one, as it would take an exponential amount
382 of computational time. Cases characterised by a high degree of non-linearity to which little information
383 is available to bound the search, such as incorporating interactions in SDMs (Kissling et al., 2012), can
384 be classified as complex problems. The ability to deal with this complexity is an advantage over other
385 methods. EAs also have notable weaknesses and pitfalls which are discussed in this section. It is important
386 to note that other potentially suitable metaheuristic algorithms exist such as particle swarm optimization
387 and ant colony optimization. They are also shortly discussed as they share a number of characteristics with
388 EAs making them interesting for machine learning in species distribution modelling.

390 4.1. Problem encoding and flexibility

391 A clear strength is the flexibility of EAs offering the chance to implement any type of machine learning
392 problem by using the encoding-model interface. Specifically, the ability of (1) encoding the model in a
393 computation element, the chromosome, and (2) using mappers to translate chromosomes to models allows
394 separating the process of training (with operators) from fitness calculation (model performance evalua-
395 tion). This flexibility has already been illustrated in GARP, where different relations, *e.g.* logistic, linear
396 or boolean, can be used in the software (Olden et al., 2008).

397
398 The mentioned flexibility allows to define various ways of model training, *i.e.* estimating model parameters
399 and/or reducing model complexity. Reducing model complexity in conjunction with learning can present an
400 opportunity for the use of individual- and agent-based to support species distribution modelling. Indeed,
401 the structure of these individual- and agent-based models can be complex (Grimm et al., 2010) and model
402 simplification with flexible machine learning algorithms could allow for a further automation of model de-
403 velopment. To reduce model complexity, the most relevant features of a model are selected by encoding
404 embedded or wrapper feature selection (Saeys et al., 2007). Wrapper feature selection is concerned with
405 selecting features for other data-driven or an already parameterized model. As opposed to this, embed-
406 ded feature selection estimates model parameters and selects features simultaneously. For wrapper feature
407 selection, a binary string encoding the inclusion (1) or exclusion (0) can be used (D’heygere et al., 2006)
408 whereas, for embedded feature selection, a ‘list of list’ approach can be used (Gobeyn, 2018). For encoding
409 a binary or continuous string, the reader is referred to Haupt and Haupt (2004). The list of list, a first order
410 list is implemented in the genotype to represent a feature in- or exclusion. If an inclusion for a feature is

411 considered then a second-order list is defined, holding the value for the model parameters coupled to the in-
412 put feature, *i.e.* coefficients of the response curve. The list of list approach seems to be promising, however,
413 additional research is required to verify its performance. For parameter estimation, a string of continuous
414 values (for example, values of model parameters describing the species' niche, see (Fukuda et al., 2011))
415 can be implemented in the genotype of the chromosomes. These are then translated to model parameters
416 and -after model execution- a fitness value.

417

418 A disadvantage of the chromosome encoding and the use of a mapper function is that a certain amount
419 of programming skills is required. This might hinder novel users to use EAs or other metaheuristics for
420 their machine learning application. However, since machine learning with EAs is specifically applicable
421 to increase transparency of complex models, it is expected that the initial investment in programming
422 will be the better option - especially in the long run. Even more, open science is challenging ecological
423 informaticians to increase code flexibility, modularity and transparency (Golding et al., 2017) leading to
424 a more user-friendly experience in programming languages such as R and Python. In addition, a number
425 of initiatives are taken in the field of computer science to help non-expert and expert users to deal with
426 feature selection and hyperparameter optimization. For example, Auto-WEKA (Feurer et al., 2015) and
427 Auto-skLearn (Kotthoff et al., 2016) are initiatives that consider the problem of simultaneously selecting
428 a learning algorithm and values for the hyperparameters through Bayesian learning. As such, these tools
429 offer the opportunity to investigate the position of EAs in comparison to other machine learning methods
430 for specific problems tackled in species distribution modelling and ecology in general.

431 *4.2. Population-based approach*

432 The population-based approach of EAs is considered the second advantage for species distribution mod-
433 elling since ecological phenomena characterised by a large amount of noise are too complex to describe by
434 one model (Fukuda et al., 2013; Merow et al., 2014; Muñoz-Mas et al., 2017; Vezza et al., 2015). Using
435 multiple models in the context of ensemble learning is useful to reflect model uncertainties (Araújo and
436 New, 2007). Practically, the EA would be run a couple of times preferably with other samples of the
437 training data (*i.e.* cross-validation or bootstrapping) and track the best models found in each run. En-
438 semble learning has shown to be valuable to avoid SDM overfitting, especially for modelling rare species
439 (Breiner et al., 2015). The population-based approach of EAs allows providing an informative ensemble of
440 near-optimal solutions rather than just one optimal solution. In this perspective, EAs can be used to gener-

441 ate ensembles comparable (Muñoz-Mas et al., 2017) and possibly serve as an alternative for the RF method.

442

443 The combination of iterating a number of solutions and the crossover and mutation operators offers the
444 opportunity to explore multiple areas of the search space (Holland, 2000). Applied to feature selection,
445 it allows tracking interesting combinations of features over several generations. This is considered a great
446 strength and a competitive alternative to stepwise selection procedures usually used in species distribution
447 modelling. In stepwise selection procedures, an alternative model is tested to data by iteratively excluding
448 (including) a feature (Zuur et al., 2009). These approaches are considered greedy because they make locally
449 optimal decisions with the assumption that a (near-)optimal solution will be found in the vicinity of this
450 local solution. Although the forward selection approach is computationally efficient, this procedure may
451 ignore informative combinations of features which are individually only marginally relevant. The search
452 behaviour of EAs is totally different: They combine and test solutions that are located in various regions
453 of the search space.

454

455 EAs and ant colony optimization are population-based approaches able to deal with combinatorial opti-
456 mization problems (Boussaïd et al., 2013) whereas particle swarm optimization was initially designed to
457 solve continuous problems (Kennedy and Eberhart, 1997). Combinatorial optimization problems are a
458 class of discrete optimization problems in which the input arguments encode permutations, combinations
459 or variations (Scheerlinck et al., 2009). The way the candidate solutions are generated is the main difference
460 between EAs and ant colony optimization. In EAs, candidate solutions are encoded as strings of bits or real
461 numbers of the chromosomes whereas for ant colony optimization the potential solutions are encoded in the
462 environment of ants. That is, the ants or agents propagate through the search space and new candidate
463 solutions are constructed from the information in this environment. This way, the memory of the system is
464 embedded in the environment rather than the objects. This property makes ant colony optimization more
465 appealing for modelling dynamically changing systems (Maier et al., 2003; Szemis et al., 2012; Zecchin
466 et al., 2006).

467

468 For now, the application of ant colony optimization in species distribution modelling might seem less in-
469 teresting since data are often not available over multiple time instances. As depicted in the introduction,
470 near-real-time data are expected to arrive as technologies in species-tracking and remote sensing are con-

471 tinuously improved (Cord et al., 2014; Pauwels et al., 2014; Bastille-Rousseau et al., 2017). As ant colony
472 optimization is able to deal with dynamic constraints without reinitialisation, it is expected to be appro-
473 priate to deal with these type of dynamic data. In these cases, the use of ant colony optimization for model
474 identification could be assessed as superior to EAs.

475

476 A fairly novel class of population-based methods are ‘Estimation of distribution algorithms’ (EDAs). These
477 algorithms guide the search for an optimal solution by sampling probabilistic models of candidate solutions,
478 and by using selection operators also applied in EAs. The aim of EDAs is not only to optimize models,
479 but also to provide a series of probabilistic models revealing characteristics of the problem being solved
480 (Pelikan et al., 2015). Other examples of newly developed population-based methods to obtain this type of
481 information are ‘irace’ (López-Ibáñez et al., 2016) and sequential model-based optimization (Hutter et al.,
482 2011). They all share the aim of automatic algorithm configuration, defined as finding good algorithm
483 settings (values for hyperparameters, operators) for solving unseen problem instances by learning on a set
484 of training problem instances (López-Ibáñez et al., 2016). Applying this type of algorithms to train SDMs
485 could be interesting to further learn about the characteristics of the training problem at hand. Sample
486 prevalence is a typical example of a characteristic of a training data set (see also Table 2). The mentioned
487 techniques could thus be used to train models on data sets with varying sample prevalence so to provide
488 interesting insights on the effect of sample prevalence on - not only the objective measure - but also algo-
489 rithms’ functioning.

490

491 *4.3. Hyperparameters*

492 The standard application of an EA requires five hyperparameters to be optimized or tuned (population
493 size, a stopping criterion, selection rate and crossover and mutation rate). This can be considered as a
494 disadvantage since the performance of the EA depends on the choice of these hyperparameters (Grefenstette,
495 1986; Feurer et al., 2015). Guidelines for (automated) tuning these hyperparameters are found in the
496 literature (Gibbs et al., 2008, 2010; López-Ibáñez et al., 2016). Yet, it is important to note that the ‘No
497 Free Lunch Theorem’ states that there is no global set of hyperparameters effective for every problem
498 (Wolpert and Macready, 1997). Consequently, every class of problems will require hyperparameter testing.
499 The results our literature review show that a limited number of studies (eight out of 25) used an iterative
500 approach to obtain hyperparameter values. In addition, no significant relation between hyperparameters

501 could be identified (see supportive information 3). This is in line with the findings of Gibbs et al. (2008)
502 who empirically determined the degree of interaction between hyperparameters for a list of optimization
503 problems. Only the population size shows a strong inverse relationship with the mutation rate whereas
504 the interaction between other hyperparameters was found not to be as relevant for the GA performance.
505 Within our analysis, we could not determine a relation between the number of chromosomes and the
506 mutation rate. The reason for this observation is that hyperparameters are rarely optimized in the field of
507 species distribution modelling. Algorithms are used, and settings seem to be copied from other publications
508 without explicit reasoning (see Table 2: Boets et al. (2013); D’heygere et al. (2003, 2006) and Zarkami et al.
509 (2012)). As such, we suspect that hyperparameters values in the studies in Table 2 are sub-optimal. Here,
510 we advocate the practice of testing and reporting the values for hyperparameters and their effect on the
511 objective function, so readers can assess which hyperparameters might be useful for a specific application in
512 species distribution modelling. We promote the use of guidelines to have a good estimate of optimal values
513 for the hyperparameters as those in Gibbs et al. (2008) and Gibbs et al. (2010). Although the number of
514 hyperparameters to be determined may be a weakness of EAs, many metaheuristic algorithms (*i.e.* ant
515 colony optimization, particle swarm optimization, simulated annealing) and machine learning algorithms
516 share this shortcoming. As noted at the end of section 4.1, a number of tools developed in computer science
517 are being developed to automate the hyperparameter optimization problem (Feurer et al., 2015).

518 4.4. Multiobjective machine learning

519 An opportunity of EAs in species distribution modelling is their potential use as multiobjective machine
520 learning methods which aim to train a model based on multiple -potentially conflicting - objectives. Typi-
521 cally, the purpose of species distribution modelling is to train models which estimate species presence and
522 absence well. In many cases, it is desired - for instance in decision management - to give a higher weight to
523 one or the other (Mouton et al., 2009). A number of evaluation criteria based on the classification of the
524 occurrence probability (*e.g.* Cohen’s Kappa or True Skill Statistics) are being used to pool the degree of
525 correct estimation of species presence and absence (Mouton et al., 2010). Unfortunately, these evaluation
526 measures depend on sample prevalence. Consequently, training models with these data having varying
527 sample prevalence are biased. A pragmatic approach to solve this issue is to keep sample prevalence equal
528 over all data samples and/or to define a trade-off between commission and omission errors in the objective
529 function (Allouche et al., 2006; Manel et al., 2001; Mouton et al., 2010).

530

531 The trade-off between omission and commission errors can be viewed as a multiobjective problem. EAs have
532 proven to be adequate techniques to identify trade-offs between objectives (Penn et al., 2013; Sweetapple
533 et al., 2014). In general, EAs can be used to determine the entire set of Pareto optimal solutions or
534 at least a representative subset. A Pareto optimal set is a set of solutions that are nondominated when
535 compared with other solutions of the solution space (Deb et al., 2000). For example, for species distribution
536 modelling, a Pareto optimal set could be a set of equally valid solutions to a problem presenting the trade-off
537 between commission and omission errors. This way, decision makers obtain a set of solutions that can be
538 very valuable for different aspects of ecosystem decision management (Guisan et al., 2013). A well-known
539 example of a multiobjective optimizer using an EA is the non-dominated sorting GA II of Deb et al. (2000).
540 In this algorithm, a simple GA with uniform crossover and mutation but with specific selection operators is
541 used. For the selection operator, different nondominant fronts are identified. These nondominant fronts are
542 estimates of the Pareto front defined by two or more objectives. The chromosomes in each non-dominant
543 front have the same assigned dummy fitness value, ranked according to the ‘strength’ of the front. These
544 dummy fitness values are used to select chromosomes (Deb et al., 2000). This process is repeated until
545 a nondominant front equal or close to the Pareto optimal front is found. An example of the use of the
546 non-dominated sorting GA II in ecology is presented by Côté et al. (2007).

547 **5. Recommendations for application**

548 EAs and other metaheuristic algorithms are particularly useful to solve problems such as feature selection,
549 parametrisation of complex models, and optimization of other learning algorithms. These algorithms are
550 likely not suited to solve every problem as the development of a specific EA will require high investment
551 costs - in terms of programming and algorithm understanding - returning little improvement in model
552 insight and predictive performance. In these cases, the use of machine learning methods, such as decision
553 trees, GLMs or Maxent would be more appropriate. However, we recommend to consider EAs when the
554 problem at hand has one of the following characteristics:

- 555 • The problem and search for a solution is expected to be complex (*e.g.* includes species interactions
556 or many features) and little information is available to *a priori* reduce complexity (see, for example,
557 Kissling et al. (2012)).
- 558 • Many (complex) boundaries *can* be formulated for the problem. These could, for instance, be obtained
559 from experts or ecological databases (Verberk et al., 2012).

- 560 • Solutions to the problem are required to be transparent and flexible for model (re-)analysis, for
561 instance for decision management (Adriaenssens et al., 2004).
- 562 • The input data set has a high number of features, and manual feature reduction is no longer possible
563 (*e.g.* pesticide database of river sediment in Flanders counts more than 200 identified pesticides
564 (VMM, 2018)).
- 565 • The model knows many parameters which have to be calibrated (Van Broekhoven et al., 2007).
- 566 • A trade-off between objective functions is required for decision management applications. This can,
567 for instance, be the trade-off between model complexity and performance, or between the correct
568 estimation of species presence and absence.

569 For a specific problem, one can select from a number of EA implementations. In Table 4, a suggestion
570 for the type of EA are provided for a number of problems. Two ‘trivial’ problems are listed, parameter
571 estimation and feature selection (see row one and two), whereas other applications are less obvious, and
572 often problem-specific implementations. For the calibration of parameter-rich models (> 10 parameters), a
573 binary or real GA encoded can be used, since both are expected to perform equally well (Van Broekhoven
574 et al., 2007). The second case involves the reduction of the number of input features with the help of
575 EAs. Typically, this applies to data sets for which a large number of potential input variables can explain
576 species occurrence. This type of learning could be particularly interesting when remote sensing products
577 are used, in order to reduce the amount of input data required to estimate species distributions from spa-
578 tial input data (Hampton et al., 2013). Automated variable selection with EAs can be helpful to steer
579 model development, but as noted by Araújo and Guisan (2006), this should not replace a selection based
580 on expert knowledge. In the case of feature selection, a binary encoded GA is implemented, encoding the
581 in- or exclusion of input features (D’heygere et al., 2006). This feature selection can be helpful for the
582 optimization of stacked SDMs or population-based SDMs. In these SDMs, models for different species are
583 coupled with each other (Guisan and Rahbek, 2011), and are allowed to interact. With the number of
584 species considered in these stacked SDMs, the number of model elements increases exponentially (due to
585 one-on-one interaction). Binary GAs can be used to simplify these models, preventing an overly complex
586 model to be fitted to a limited number of species occurrence observations. In addition, binary GAs can be
587 used to optimize artificial neural networks. In this case, different layers or neurons can be implemented in
588 the binary string, and the structure of the ANN can be optimized (see Muñoz-Mas et al. (2017)).

589

590 For simple binary and real-coded GAs, a selection rate of 0.5, a crossover rate above 0.8, mutation rate
591 lower than 0.2 and 100 generations will in general work well, when the number of chromosomes is between
592 30 and 200, independent of the chromosome length (Gibbs et al. (2008) and Table 2). For the choice of
593 the selection operator, we advise using tournament or roulette wheel selection. Both are simple to under-
594 stand, and generally give satisfying results when compared to other selection operators (Goldberg and Deb,
595 1991). The use of elitism is advised, however, it is important to note that the use of elitism can decrease
596 the population diversity, and facilitate faster convergence. The need for fast convergence, motivated by
597 limited available computation resources can be an important boundary condition in choosing the number
598 of model evaluations. This number is determined by the number of chromosomes multiplied by a number of
599 generations. A limited number of model evaluations, 400 and 2500, have been used, and have presumably
600 led to satisfying results. Increasing the number of evaluations can be usefull, however, it is possible that
601 gains in accuracy or precision are marginal. As discussed by Gibbs et al. (2008), the number of evaluations
602 should vary as a function of the available computational resources. As a rule of thumb, we advise to focus
603 on a cross-validation resampling strategy and on repeated execution of the EA/cross-validation strategy to
604 increase robustness, rather than employing a larger number of model evaluations.

605

606 Users are advised to consider stratification according to sample prevalence to a design cross validation
607 strategy. As discussed above, accuracy measures can vary as a function of this prevalence. To make results
608 comparable, it is of importance that data are stratified according to this prevalence. The choice for a
609 number of folds, and repetitions will depend on the available computation resources, and the size of the
610 data set. Precision will increase with a higher number of repetitions and folds, leading to a longer runtime.
611 When data sets are small, and models are learned fast, a higher number is thus preferred. In contrast, when
612 learning is slow, one can opt to choose fewer folds and repetitions (Kohavi, 1995). As discussed above, one
613 can also consider lowering the number of function evaluations.

614

615 When further fine-tuning of the hyperparameters is desired, we recommend using the guidelines of Gibbs
616 et al. (2008), as every specific problem can have a unique set of optimal hyperparameters. The choice of
617 the objective function, which the GA has to optimize, depends on the study objectives: does one aim to
618 estimate species presence well, or rather absence? If the former is true than a higher weight should be given
619 to sensitivity, in the case of the latter, specificity (see Table 3). In case a trade-off between both should

620 be identified, one can consider a multi-objective EA. In these algorithms, a specific selection operator is
621 implemented in the GA to weight different objectives (see for example the non-dominated sort in NSGA-II
622 (Deb et al., 2002)).

623

624 Three main points need to be taken into account when machine learning or other algorithms are consid-
625 ered to solve a hypothesis. First, specifying the model and its structural component, the subject of model
626 training and the objective of the model and the study (see Box 3) is important (Guisan and Zimmermann,
627 2000). For example, is the aim of the model to understand a specific theoretical assumption about species
628 interaction? Or is the aim to develop a predictive model for estimating species occurrence in an ecosystem
629 with many interactions? In a second step, an algorithm to train the model(s) needs to be selected (Box 3,
630 second part). Specifically, algorithm operators, problem encoding and operators need to be defined. Here,
631 it is important to make a distinction between algorithms which make use of explicit encoding (EAs, ant
632 colony optimization) and those which do not (decision trees, GLMs). It is empirically found that algorithms
633 making use of encoding work well to train models with hypotheses embedded in the model structure (Maier
634 et al., 2014). The initial choice for a type of algorithm and use of encoding will hence determine the choice
635 for hyperparameters and operators.

636

637 Finally, a platform to implement the approach for the machine learning application is required. GUI
638 packages can be used, however, adopting these packages can considerably limit the options which make
639 EAs interesting in the first place. For that reason, we advise to use a high-level scripting language such as
640 Python or R and search for existing codes implemented in these languages. An additional advantage of using
641 high-level scripting languages for machine learning applications is their transferability to high performance
642 and cloud computing infrastructure. Preferably, the scripting is done in an open science context allowing
643 for continuous code improvement and validation through modular scripting. For a good introduction on
644 modular scripting applied to ecology, we refer to Golding et al. (2017). Typically, open science is performed
645 on code sharing platforms such as GitHub (<https://github.com/>).

646 **6. Future perspectives and conclusions**

647 Recent advances in theoretical ecology (Leibold et al., 2004) and conceptual modelling (Guisan and Rahbek,
648 2011) are challenging scientists to continuously develop new ways to deal with this increasing complexity.

649 The field of machine learning has proven to be useful to tackle these questions, despite that researchers are
650 struggling to identify the appropriate approach to address increasing complexity (Kissling et al., 2012).
651 Maxent is currently the most used technique to model species distributions when considering terrestrial
652 cases. For freshwater system case studies, innovative methods such as artificial neural networks and EAs
653 are increasingly being used to solve less straightforward problems. Model developers will be required to
654 deal with this increased complexity, preferably in an open science context. This depicts full transparency
655 in the methodology, but also the practical encoding (Golding et al., 2017; Phillips et al., 2017). In addition,
656 it requires the developed algorithms to be easily transferable and adaptable to new problems. Considering
657 these aspects, EAs and other metaheuristic algorithms are found to be of particular use since they split
658 the training process from the objective function evaluation (model run). In addition, EAs allow dealing
659 with complex cases (Eiben and Smith, 2015) making them appropriate candidates to train the next gen-
660 eration of species distribution models. Dealing with hyperparameters optimization and the requirement
661 of programming and modelling skills are considered disadvantageous, hindering the use of EAs and other
662 metaheuristic algorithms. For the first, hyperparameter optimization methods are already available giving
663 satisfying results for multiple problems (Gibbs et al., 2008; Gobeyn et al., 2017). The second, the need for
664 modelling know-how will require standardization, documentation and refinement of the algorithm devel-
665 opment and application process (Jakeman et al., 2006; Grimm et al., 2010) going hand in hand with the
666 philosophy of open science. In this review paper, a number of suggestions with respect to the definition
667 of the model, algorithm and implementation are given. With this, we aim to stimulate ecologists to use
668 and further refine the development of EAs applied to species distribution modelling, hypothesis testing and
669 preferably ecology in general.

670

671 As technological advances in machine learning are reshaping the way scientist develop models and analyse
672 data, researches are increasingly aware that one specific algorithm won't offer a tailor-made solution to
673 every problem (Chatfield, 1993). With this synthesis, an insight is presented on how to use EAs as a
674 technique to solve specific problems in ecology rather than using it as a ready-to-use technology to map
675 species distributions.

676 **Acknowledgments**

677 The authors would like to thank two anonymous reviewers for their thorough review and valuable comments
678 which have considerably improved the manuscript. This study was carried out within the CROSSLINK
679 project, funded through BiodivERsA, under the Horizon 2020 ERA-NET COFUND scheme. This research
680 was partly funded through the 2013/2014 BiodivERsA/FACCE-JPI joint call, with the national funder
681 BMBF - German Federal Ministry of Education and Research (Project TALE Towards multifunctional
682 agricultural landscapes in Europe: Assessing and governing synergies between food production, biodiversity,
683 and ecosystem services, grant 01 LC 1404 A).

684 **Author contribution**

685 S.G., A.M.M., A.F.C., M.V. and P.L.M.G. designed the research. S.G. conducted the literature review. S.G.
686 wrote the manuscript, and A.M.M., A.F.C., A.K., M.V. and P.L.M.G. provided edits to the manuscript.

687 **References**

- 688 Adriaenssens, V., 2004. Knowledge-based macroinvertebrate habitat suitability models for use in ecological
689 river management. Ph.D. thesis. Ghent University.
- 690 Adriaenssens, V., De Baets, B., Goethals, P.L.M., De Pauw, N., 2004. Fuzzy rule-based models for decision
691 support in ecosystem management. *Science of the Total Environment* 319, 1–12.
- 692 Allouche, O., Tsoar, A., Kadmon, R., 2006. Assessing the accuracy of species distribution models: Preva-
693 lence, kappa and the true skill statistic (TSS). *Journal of Applied Ecology* 43, 1223–1232.
- 694 Araújo, M.B., Guisan, A., 2006. Five (or so) challenges for species distribution modelling. *Journal of*
695 *Biogeography* 33, 1677–1688.
- 696 Araújo, M.B., New, M., 2007. Ensemble forecasting of species distributions. *Trends in Ecology and*
697 *Evolution* 22, 42–47.
- 698 Austin, M.P., 2007. Species distribution models and ecological theory: A critical assessment and some
699 possible new approaches. *Ecological Modelling* 200, 1–19.

- 700 Austin, M.P., Cunningham, R.B., 1981. Observational analysis of environmental gradients. *Proceedings of*
701 *the Ecological Society of Australia* 11, 109–119.
- 702 Baert, J.M., Janssen, C.R., Sabbe, K., De Laender, F., 2016. Per capita interactions and stress tolerance
703 drive stress-induced changes in biodiversity effects on ecosystem functions. *Nature Communications* 7,
704 1–8.
- 705 Bastille-Rousseau, G., Murray, D.L., Schaefer, J.A., Lewis, M.A., Mahoney, S., Potts, J.R., 2017. Spatial
706 scales of habitat selection decisions: implications for telemetry-based movement modelling. *Ecography*
707 40, 1–7.
- 708 Bennetsen, E., Gobeyn, S., Goethals, P.L.M., 2016. Species distribution models grounded in ecological
709 theory for decision support in river management. *Ecological Modelling* 325, 1–12.
- 710 Bergstra, J., Bengio, Y., 2012. Random Search for Hyper-Parameter Optimization. *Journal of Machine*
711 *Learning Research* 13, 281–305.
- 712 Boets, P., Holguin, G., Lock, K., Goethals, P.L.M., 2013. Data-driven habitat analysis of the Ponto-Caspian
713 amphipod *Dikerogammarus villosus* in two invaded regions in Europe. *Ecological Informatics* 17, 36–45.
- 714 Booth, T.H., 1985. A new method to assist species selection. *The Commonwealth Forestry Review* 64,
715 241–250.
- 716 Booth, T.H., Nix, H.A., Busby, J.R., Hutchinson, M.F., 2014. Bioclim: The first species distribution
717 modelling package, its early applications and relevance to most current MaxEnt studies. *Diversity and*
718 *Distributions* 20, 1–9.
- 719 Boussaïd, I., Lepagnot, J., Siarry, P., 2013. A survey on optimization metaheuristics. *Information Sciences*
720 237, 82–117.
- 721 Breiman, L., 2001. Random Forests. *Machine Learning* 45, 5–32.
- 722 Breiman, L., Friedman, J.H., Olshen, R.A., Stone, C.J., 1984. *Classification and regression trees*. 1 ed.,
723 Taylor & Francis.
- 724 Breiner, F.T., Guisan, A., Bergamini, A., Nobis, M.P., 2015. Overcoming limitations of modelling rare
725 species by using ensembles of small models. *Methods in Ecology and Evolution* 6, 1210–1218.

- 726 Carpenter, G., Gillison, A.N., Winter, J., 1993. DOMAIN: a flexible modelling procedure for mapping
727 potential distributions of plants and animals. *Biodiversity and Conservation* 2, 667–680.
- 728 Chatfield, C., 1993. Neural networks: Forecasting breakthrough or passing fad? *International Journal of*
729 *Forecasting* 9, 1–3.
- 730 Chen, W.C., Chang, N.B., Chen, J.C., 2003. Rough set-based hybrid fuzzy-neural controller design for
731 industrial wastewater treatment. *Water Research* 37, 95–107.
- 732 Cord, A.F., Klein, D., Gernandt, D.S., de la Rosa, J.A.P., Dech, S., 2014. Remote sensing data can improve
733 predictions of species richness by stacked species distribution models: A case study for Mexican pines.
734 *Journal of Biogeography* 41, 736–748.
- 735 Côté, P., Parrott, L., Sabourin, R., 2007. Multi-objective optimization of an ecological assembly model.
736 *Ecological Informatics* 2, 23–31.
- 737 Cutler, D., Edwards, T., Beard, K.H., Cutler, A., Hess, K., Gibson, J., 2007. Random Forests for Classifi-
738 cation in Ecology. *Ecology* 88, 2783–2792.
- 739 D’Angelo, D.J., Howard, L.M., Meyer, J.L., Gregory, S.V., Ashkenas, L.R., 1995. Ecological uses for genetic
740 algorithms: predicting fish distributions in complex physical habitats. *Canadian Journal of Fisheries and*
741 *Aquatic Sciences* 52, 1893–1908.
- 742 DeAngelis, D.L., Yurek, S., 2017. Spatially Explicit Modeling in Ecology: A Review. *Ecosystems* 20,
743 284–300.
- 744 Deb, K., Agrawal, S., Pratap, A., Meyarivan, T., 2000. A Fast Elitist Non-dominated Sorting Genetic
745 Algorithm for Multi-objective Optimization: NSGA-II. Springer Berlin Heidelberg, Berlin, Heidelberg.
746 pp. 849–858.
- 747 Deb, K., Pratap, A., Agarwal, S., Meyarivan, T., 2002. A fast and elitist multiobjective genetic algorithm:
748 NSGA-II. *IEEE Transactions on Evolutionary Computation* 6, 182–197.
- 749 D’heygere, T., Goethals, P.L.M., De Pauw, N., 2003. Use of genetic algorithms to select input variables in
750 decision tree models for the prediction of benthic macroinvertebrates. *Ecological Modelling* 160, 291–300.

751 D'heygere, T., Goethals, P.L.M., De Pauw, N., 2006. Genetic algorithms for optimisation of predictive
752 ecosystems models based on decision trees and neural networks. *Ecological Modelling* 195, 20–29.

753 Ding, S., Li, H., Su, C., Yu, J., Jin, F., 2013. Evolutionary artificial neural networks: A review. *Artificial*
754 *Intelligence Review* 39, 251–260.

755 Dormann, C.F., Schymanski, S.J., Cabral, J., Chuine, I., Graham, C., Hartig, F., Kearney, M., Morin, X.,
756 Römermann, C., Schröder, B., Singer, A., 2012. Correlation and process in species distribution models:
757 bridging a dichotomy. *Journal of Biogeography* 39, 2119–2131.

758 Dunning, J.B., Stewart, D.J., Danielson, B.J., Noon, B.R., Root, T.L., Lamberson, R.H., Stevens, E.E.,
759 Danielson, B.J., Noon, B.R., Root, T.L., Lamberson, R.H., Stevens, E.E., 1995. Spatially Explicit
760 Population Models : Current Forms and Future Uses. *Ecological Applications* 5, 3–11.

761 Eiben, A.E., Smith, J., 2015. From evolutionary computation to the evolution of things. *Nature* 521,
762 476–482.

763 Elith, J., Graham, C., Anderson, R., Dudík, M., Ferrier, S., Guisan, A., Hijmans, R., Huettmann, F.,
764 Leathwick, J., Lehmann, A., Li, J., Lohmann, L., Loiselle, B., Manion, G., Moritz, C., Nakamura, M.,
765 Nakazawa, Y., Overton, J., Peterson, A., Phillips, S., Richardson, K., Scachetti-Pereira, R., Schapire,
766 R., Soberon, J., Williams, S., Wisz, M., Zimmermann, N., 2006. Novel methods improve prediction of
767 species' distributions from occurrence data. *Ecography* 29, 129–151.

768 Elith, J., Leathwick, J.R., 2009. Species Distribution Models: Ecological Explanation and Prediction Across
769 Space and Time. *Annual Review of Ecology, Evolution, and Systematics* 40, 677–697.

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

772 Everaert, G., Pauwels, I., Bennetsen, E., Goethals, P.L.M., 2016. Development and selection of decision
773 trees for water management: Impact of data preprocessing, algorithms and settings. *AI Communications*
774 29, 711–723.

775 Favaro, L., Tirelli, T., Pessani, D., 2011. Modelling habitat requirements of white-clawed crayfish (*Aus-*
776 *tropotomobius pallipes*) using support vector machines. *Knowledge and Management of Aquatic Ecosys-*
777 *tems* 401, 21.

- 778 Feurer, M., Klein, A., Eggenberger, K., Springenberg, J., Blum, M., Hutter, F., 2015. Efficient and Robust
779 Automated Machine Learning. *Advances in Neural Information Processing Systems* 28 , 2944–2952.
- 780 Fukuda, S., 2009. Consideration of fuzziness: Is it necessary in modelling fish habitat preference of Japanese
781 medaka (*Oryzias latipes*)? *Ecological Modelling* 220, 2877–2884.
- 782 Fukuda, S., De Baets, B., Mouton, A.M., Waegeman, W., Nakajima, J., Mukai, T., Hiramatsu, K., Onikura,
783 N., 2011. Effect of model formulation on the optimization of a genetic Takagi-Sugeno fuzzy system for
784 fish habitat suitability evaluation. *Ecological Modelling* 222, 1401–1413.
- 785 Fukuda, S., De Baets, B., Waegeman, W., Verwaeren, J., Mouton, A.M., 2013. Habitat prediction and
786 knowledge extraction for spawning European grayling (*Thymallus thymallus* L.) using a broad range of
787 species distribution models. *Environmental Modelling & Software* 47, 1–6.
- 788 Fukuda, S., Hiramatsu, K., 2008. Prediction ability and sensitivity of artificial intelligence-based habitat
789 preference models for predicting spatial distribution of Japanese medaka (*Oryzias latipes*). *Ecological*
790 *Modelling* 215, 301–313.
- 791 Fukuda, S., Mouton, A.M., De Baets, B., 2012. Abundance versus presence/absence data for modelling fish
792 habitat preference with a genetic Takagi-Sugeno fuzzy system. *Environmental Monitoring and Assessment*
793 184, 6159–6171.
- 794 Gendreau, M., Potvin, J.Y., 2010. *Handbook of Metaheuristics*. 2 ed., Springer International Publishing,
795 New York.
- 796 Gibbs, M.S., Dandy, G.C., Maier, H.R., 2008. A genetic algorithm calibration method based on convergence
797 due to genetic drift. *Information Sciences* 178, 2857–2869.
- 798 Gibbs, M.S., Maier, H.R., Dandy, G.C., 2010. Comparison of Genetic Algorithm Parameter Setting Methods
799 for Chlorine Injection Optimization. *Journal of Water Resources Planning and Management* 136, 288–
800 291.
- 801 Gibbs, M.S., Maier, H.R., Dandy, G.C., 2015. Using characteristics of the optimisation problem to deter-
802 mine the Genetic Algorithm population size when the number of evaluations is limited. *Environmental*
803 *Modelling & Software* 69, 226–239.

804 Gobeyn, S., 2018. Species Distribution Model Identification Tool. URL: [https://github.com/](https://github.com/Sachagobeyn/SDMIT)
805 Sachagobeyn/SDMIT.

806 Gobeyn, S., Goethals, P.L., 2017. A variable length chromosome genetic algorithm approach to identify
807 species distribution models useful for freshwater ecosystem management, in: Denzer, R., Schimak, G.,
808 Hebíček, J. (Eds.), Environmental Software Systems. Infrastructures, Services and Applications. Springer
809 International Publishing, Cham, pp. 196–208.

810 Gobeyn, S., Volk, M., Dominguez-Granda, L., Goethals, P.L.M., 2017. Input variable selection with a
811 simple genetic algorithm for conceptual species distribution models: A case study of river pollution in
812 Ecuador. *Environmental Modelling & Software* 92, 269–316.

813 Goethals, P.L.M., Dedecker, A.P., Gabriels, W., Lek, S., De Pauw, N., 2007. Applications of artificial
814 neural networks predicting macroinvertebrates in freshwaters. *Aquatic Ecology* 41, 491–508.

815 Goldberg, D.E., Deb, K., 1991. A Comparative Analysis of Selection Schemes Used in Genetic Algorithms.
816 *Foundations of Genetic Algorithms* 1, 69–93.

817 Golding, N., August, T.A., Lucas, T.C., Gavaghan, D.J., van Loon, E.E., Mcinerny, G., 2017. The zoon r
818 package for reproducible and shareable species distribution modelling. *Methods in Ecology and Evolution*
819 9, 1–9.

820 Grefenstette, J.J., 1986. Optimization of Control Parameters for Genetic Algorithms. *IEEE Transactions*
821 *on Systems, Man and Cybernetics* 16, 122–128.

822 Grimm, V., Berger, U., DeAngelis, D.L., Polhill, J.G., Giske, J., Railsback, S.F., 2010. The ODD protocol:
823 A review and first update. *Ecological Modelling* 221, 2760–2768.

824 Guisan, A., Rahbek, C., 2011. SESAM a new framework integrating macroecological and species distribu-
825 tion models for predicting spatio-temporal patterns of species assemblages. *Journal of Biogeography* 38,
826 1433–1444.

827 Guisan, A., Thuiller, W., 2005. Predicting species distribution: Offering more than simple habitat models.
828 *Ecology Letters* 8, 993–1009.

829 Guisan, A., Tingley, R., Baumgartner, J.B., Naujokaitis-Lewis, I., Sutcliffe, P.R., Tulloch, A.I.T., Regan,
830 T.J., Brotons, L., McDonald-Madden, E., Mantyka-Pringle, C., Martin, T.G., Rhodes, J.R., Maggini, R.,

831 Setterfield, S.A., Elith, J., Schwartz, M.W., Wintle, B.A., Broennimann, O., Austin, M.P., Ferrier, S.,
832 Kearney, M.R., Possingham, H.P., Buckley, Y.M., 2013. Predicting species distributions for conservation
833 decisions. *Ecology Letters* 16, 1424–1435.

834 Guisan, A., Zimmermann, N.E., 2000. Predictive habitat distribution models in ecology. *Ecological Mod-*
835 *elling* 135, 147–186.

836 Hampton, S.E., Strasser, C.A., Tewksbury, J.J., Gram, W.K., Budden, A.E., Batcheller, A.L., Duke, C.S.,
837 Porter, J.H., 2013. Big data and the future of ecology. *Frontiers in Ecology and the Environment* 11,
838 156–162.

839 Haupt, R.L., Haupt, S.E., 2004. *Algorithms Practical Genetic Algorithms*. 2 ed., John Wiley & Sons, Inc.,
840 Hoboken.

841 He, K.S., Bradley, B.A., Cord, A.F., Rocchini, D., Tuanmu, M.N., Schmidlein, S., Turner, W., Wegmann,
842 M., Pettorelli, N., 2015. Will remote sensing shape the next generation of species distribution models?
843 *Remote Sensing in Ecology and Conservation* 1, 4–18.

844 Hirzel, A.H., Le Lay, G., 2008. Habitat suitability modelling and niche theory. *Journal of Applied Ecology*
845 45, 1372–1381.

846 Hoang, T.H., Lock, K., Mouton, A.M., Goethals, P.L.M., 2010. Application of classification trees and
847 support vector machines to model the presence of macroinvertebrates in rivers in Vietnam. *Ecological*
848 *Informatics* 5, 140–146.

849 Holland, J.H., 2000. Building blocks, cohort genetic algorithms, and hyperplane-defined functions. *Evolu-*
850 *tionary Computation* 8, 373–91.

851 Hutchinson, E.G., 1957. Concluding remarks. *Cold Spring Harbor Symposia on Quantative Biology* 159,
852 415–427.

853 Hutter, F., Hoos, H.H., Leyton-Brown, K., 2011. Sequential model-based optimization for general algo-
854 rithm configuration, in: *Proceedings of the 5th International Conference on Learning and Intelligent*
855 *Optimization*, Springer-Verlag, Berlin, Heidelberg. pp. 507–523.

856 Iverson, L.R., Prasad, A.M., 1998. Predicting abundance of 80 tree species following climate change in the
857 eastern United States. *Ecological Monographs* 68, 465–485.

- 858 Jakeman, A.J., Letcher, R.A., Norton, J.P., 2006. Ten iterative steps in development and evaluation of
859 environmental models. *Environmental Modelling and Software* 21, 602–614.
- 860 Jeong, K.S., Jang, J.D., Kim, D.K., Joo, G.J., 2011. Waterfowls habitat modeling: Simulation of nest site
861 selection for the migratory Little Tern (*Sterna albifrons*) in the Nakdong estuary. *Ecological Modelling*
862 222, 3149–3156.
- 863 Kacprzyk, J., Pedrycz, W., 2015. *Springer Handbook of Computational Intelligence*. 1 ed., Springer-Verlag,
864 Berlin, Heidelberg.
- 865 Kennedy, J., Eberhart, R.C., 1997. A discrete binary version of the particle swarm algorithm, in: 1997 IEEE
866 International Conference on Systems, Man, and Cybernetics. *Computational Cybernetics and Simulation*,
867 pp. 4104–4108.
- 868 Kissling, W.D., Dormann, C.F., Groeneveld, J., Hickler, T., Kühn, I., McInerny, G.J., Montoya, J.M.,
869 Römermann, C., Schiffers, K., Schurr, F.M., Singer, A., Svenning, J.C., Zimmermann, N.E., O’Hara,
870 R.B., 2012. Towards novel approaches to modelling biotic interactions in multispecies assemblages at
871 large spatial extents. *Journal of Biogeography* 39, 2163–2178.
- 872 Kohavi, R., 1995. A study of cross-validation and bootstrap for accuracy estimation and model selection,
873 in: *International Joint Conference on Artificial Intelligence (IJCAI)*, Vol. 2, Montreal. pp. 1137–1143.
- 874 Kotthoff, L., Thornton, C., Hoos, H.H., Hutter, F., Leyton-Brown, K., 2016. Auto-WEKA 2.0: Automatic
875 model selection and hyperparameter optimization in WEKA. *Journal of Machine Learning Research* 17,
876 1–5.
- 877 Lawler, J.J., White, D., Neilson, R.P., Blaustein, A.R., 2006. Predicting climate-induced range shifts:
878 Model differences and model reliability. *Global Change Biology* 12, 1568–1584.
- 879 Leibold, M.A., Holyoak, M., Mouquet, N., Amarasekare, P., Chase, J.M., Hoopes, M.F., Holt, R.D., Shurin,
880 J.B., Law, R., Tilman, D., Loreau, M., Gonzalez, A., 2004. The metacommunity concept: A framework
881 for multi-scale community ecology. *Ecology Letters* 7, 601–613.
- 882 López-Ibáñez, M., Dubois-Lacoste, J., Pérez Cáceres, L., Birattari, M., Stützle, T., 2016. The irace package:
883 Iterated racing for automatic algorithm configuration. *Operations Research Perspectives* 3, 43–58.

- 884 Maier, H., Simpson, A., Zecchin, A., Foong, W., Phang, K., Seah, H., Tan, C., 2003. Ant Colony Optimiza-
885 tion for Design of Water Distribution Systems. *Journal of Water Resources Planning and Management*
886 129, 200–209.
- 887 Maier, H.R., Kapelan, Z., Kasprzyk, J., Kollat, J., Matott, L.S., Cunha, M.C., Dandy, G.C., Gibbs, M.S.,
888 Keedwell, E., Marchi, A., Ostfeld, A., Savic, D., Solomatine, D.P., Vrugt, J.A., Zecchin, A.C., Minsker,
889 B.S., Barbour, E.J., Kuczera, G., Pasha, F., Castelletti, A., Giuliani, M., Reed, P.M., 2014. Evolutionary
890 algorithms and other metaheuristics in water resources: Current status, research challenges and future
891 directions. *Environmental Modelling & Software* 62, 271–299.
- 892 Manel, S., Ceri Williams, H., Ormerod, S.J., 2001. Evaluating presence-absence models in ecology: The
893 need to account for prevalence. *Journal of Applied Ecology* 38, 921–931.
- 894 May, R.J., Dandy, G.C., Maier, H.R., Nixon, J.B., 2008. Application of partial mutual information variable
895 selection to ANN forecasting of water quality in water distribution systems. *Environmental Modelling &*
896 *Software* 23, 1289–1299.
- 897 McClean, C.J., Lovett, J.C., Küper, W., Hannah, L., Sommer, H.J., Wilhelm, B., Termansen, M., Smith,
898 G.F., Tokumine, S., Taplin, J.R.D., 2005. African Plant Diversity and Climate Change. *Annals of The*
899 *Missouri Botanical Garden* 92, 139–152.
- 900 McDermott, J., O'Reilly, U.m., 2015. Genetic Programming, in: Kacprzyk, J., Pedrycz, W. (Eds.), Springer
901 *Handbook of Computational Intelligence*. 1 ed.. Springer-Verlag, Berlin, Heidelberg, pp. 845–869.
- 902 McKay, R.I., 2001. Variants of genetic programming for species distribution modelling - Fitness sharing,
903 partial functions, population evaluation. *Ecological Modelling* 146, 231–241.
- 904 Merow, C., Smith, M.J., Edwards, T.C., Guisan, A., McMahon, S.M., Normand, S., Thuiller, W., Wüest,
905 R.O., Zimmermann, N.E., Elith, J., 2014. What do we gain from simplicity versus complexity in species
906 distribution models? *Ecography* 37, 1267–1281.
- 907 Meyers, G., Kapelan, Z., Keedwell, E., 2017. Short-term forecasting of turbidity in trunk main networks.
908 *Water Research* 124, 67–76.

- 909 Mount, N., Maier, H., Toth, E., Elshorbagy, A., Solomatine, D., Chang, F.J., Abrahart, R., 2016. Data-
910 driven modelling approaches for socio-hydrology: Opportunities and challenges within the Panta Rhei
911 Science Plan. *Hydrological Sciences Journal* 61, 1192–1208.
- 912 Mouton, A.M., Alcaraz-Hernández, J.D., De Baets, B., Goethals, P.L.M., Martínez-Capel, F., 2011. Data-
913 driven fuzzy habitat suitability models for brown trout in Spanish Mediterranean rivers. *Environmental*
914 *Modelling & Software* 26, 615–622.
- 915 Mouton, A.M., De Baets, B., Goethals, P.L.M., 2010. Ecological relevance of performance criteria for
916 species distribution models. *Ecological Modelling* 221, 1995–2002.
- 917 Mouton, A.M., De Baets, B., Van Broekhoven, E., Goethals, P.L.M., 2009. Prevalence-adjusted optimisation
918 of fuzzy models for species distribution. *Ecological Modelling* 220, 1776–1786.
- 919 Muñoz-Mas, R., Fukuda, S., Vezza, P., Martínez-Capel, F., 2016a. Comparing four methods for decision-
920 tree induction: A case study on the invasive Iberian gudgeon (*Gobio lozanoi*; Doadrio and Madeira,
921 2004). *Ecological Informatics* 34, 22–34.
- 922 Muñoz-Mas, R., Martínez-Capel, F., Alcaraz-Hernández, J.D., Mouton, A.M., 2017. On species distribution
923 modelling, spatial scales and environmental flow assessment with Multi-Layer Perceptron Ensembles: A
924 case study on the redfin barbel (*Barbus haasi*; Mertens, 1925). *Limnologica* 62, 161–172.
- 925 Muñoz-Mas, R., Vezza, P., Alcaraz-Hernández, J.D., Martínez-Capel, F., 2016b. Risk of invasion predicted
926 with support vector machines: A case study on northern pike (*Esox Lucius*, L.) and bleak (*Alburnus*
927 *alburnus*, L.). *Ecological Modelling* 342, 123–134.
- 928 Nelder, J.A., Wedderburn, R.W.M., 1972. Generalized Linear Models. *Journal of the Royal Statistical*
929 *Society. Series A (General)* 135, 370–384.
- 930 Nix, H.A., 1986. A biogeographic analysis of Australian elapid snakes, in: Longmore, R. (Ed.), *Australian*
931 *Flora and Fauna*. Australian Government Publishing Service, Canberra. volume 7, pp. 4–15.
- 932 Nix, H.A., McMahon, J., Mackenzie, D., 1977. No Potential areas of production and the future of pigeon pea
933 and other grain legumes in Australia, in: Wallis, E., Whiteman, P. (Eds.), *The potential for pigeon pea in*
934 *Australia: Proceedings of Pigeon Pea (*Cajanus cajan* (L.) Millsp.) Field Day*. University of Queensland,
935 Queensland. chapter 5, pp. 1–12.

- 936 Olden, J.D., Lawler, J.J.L., Poff, N.L., 2008. Machine Learning methods Without Tears: A Primer for
937 Ecologists. *The Quarterly Review of Biology* 83, 171–193.
- 938 Pal, S.K., Bandyopadhyay, S., Ray, S.S., 2006. Evolutionary computation in bioinformatics: A review.
939 *IEEE Transactions on Systems, Man, and Cybernetics Part A:Systems and Humans* 36, 601–615.
- 940 Pauwels, I.S., Goethals, P.L.M., Coeck, J., Mouton, A.M., 2014. Movement patterns of adult pike (*Esox*
941 *lucius* L.) in a Belgian lowland river. *Ecology of Freshwater Fish* 23, 373–382.
- 942 Pelikan, M., Hauschild, M.W., Lobo, F.G., 2015. Estimation of distribution algorithms, in: Kacprzyk, J.,
943 Pedrycz, W. (Eds.), *Springer Handbook of Computational Intelligence*. 1 ed.. Springer-Verlag, Berlin,
944 Heidelberg, pp. 899–928.
- 945 Penn, R., Friedler, E., Ostfeld, A., 2013. Multi-objective evolutionary optimization for greywater reuse in
946 municipal sewer systems. *Water Research* 47, 5911–5920.
- 947 Peterson, A.T., Ortega-Huerta, M.A., Bartley, J., Sanchez-Cordero, V., Soberon, J., Buddenmeier, R.H.,
948 Stockwell, D.R.B., Sánchez-Cordero, V., Soberón, J., Buddemeier, R.H., Stockwell, D.R.B., 2002. Future
949 projections for Mexican faunas under global climate change scenarios. *Nature* 416, 626–629.
- 950 Peterson, A.T., Papes, M., Eaton, M., 2007. Transferability and model evaluation in ecological niche
951 modeling: a comparison of GARP and Maxent. *Ecography* 30, 550–560.
- 952 Phillips, S.J., Anderson, R.P., Dudík, M., Schapire, R.E., Blair, M.E., 2017. Opening the black box: an
953 open-source release of Maxent. *Ecography* 40, 887–893.
- 954 Phillips, S.J., Anderson, R.P., Schapire, R.E., 2006. Maximum entropy modeling of species geographic
955 distributions. *Ecological Modelling* 190, 231–252.
- 956 Prasad, A.M., Iverson, L.R., Liaw, A., 2006. Newer classification and regression tree techniques: Bagging
957 and random forests for ecological prediction. *Ecosystems* 9, 181–199.
- 958 Prendergast, H., Hattersley, P., 1985. Distribution and cytology of Australian Neurachne and its allies
959 (*Poaceae*), a group containing C3, C4 and C3-C4 intermediate species. *Australian Journal of Botany* 33,
960 317–336.
- 961 Quinlan Ross, J., 1993. *C4. 5: Programs For Machine Learning*. doi:10.1016/S0019-9958(62)90649-6.

- 962 Rauch, W., Harremoës, P., 1999. Genetic algorithms in real time control applied to minimize transient
963 pollution from urban wastewater systems. *Water Research* 33, 1265–1277.
- 964 Rokach, L., Maimon, O., 2015. *Data Mining with Decision Trees: Theory and Applications*. 2 ed., World
965 Scientific Publishing Co. Pte. Ltd., Singapore.
- 966 Rowe, J.E., 2015. Genetic algorithms, in: Kacprzyk, J., Pedrycz, W. (Eds.), *Springer Handbook of Com-*
967 *putational Intelligence*. 1 ed.. Springer-Verlag, Berlin, Heidelberg, pp. 825–844.
- 968 Sadeghi, R., Zarkami, R., Van Damme, P., 2014. Modelling habitat preference of an alien aquatic fern,
969 *Azolla filiculoides* (Lam.), in Anzali wetland (Iran) using data-driven methods. *Ecological Modelling* 284,
970 1–9.
- 971 Sadeghia, R., Zarkami, R., Sabetraftar, K., Van Damme, P., 2013. Application of genetic algorithm and
972 greedy stepwise to select input variables in classification tree models for the prediction of habitat require-
973 ments of *Azolla filiculoides* (Lam.) in Anzali wetland, Iran. *Ecological Modelling* 251, 44–53.
- 974 Saey, Y., Inza, I., Larrañaga, P., 2007. A review of feature selection techniques in bioinformatics. *Bioin-*
975 *formatics* 23, 2507–2517.
- 976 Scheerlinck, K., Pauwels, V.R.N., Vernieuwe, H., De Baets, B., 2009. Calibration of a water and energy
977 balance model: Recursive parameter estimation versus particle swarm optimization. *Water Resources*
978 *Research* 45, W10422.
- 979 Sirbu, A., Ruskin, H., Crane, M., 2010. Comparison of evolutionary algorithms in gene regulatory network
980 model inference. *BMC Bioinformatics* 11, 59.
- 981 Srikanth, R., George, R., Warsi, N., Prabhu, D., Petry, F.E., Buckles, B.P., 1995. A variable-length genetic
982 algorithm for clustering and classification. *Pattern Recognition Letters* 16, 789–800.
- 983 Stockwell, D.R.B., Noble, I.R., 1992. Induction of sets of rules from animal distribution data: A robust
984 and informative method of data analysis. *Mathematics and Computers in Simulation* 33, 385–390.
- 985 Sutherst, R.W., Maywald, G.F., 1985. A computerised system for matching climates in ecology. *Agriculture,*
986 *Ecosystems and Environment* 13, 281–299.

- 987 Sweetapple, C., Fu, G., Butler, D., 2014. Multi-objective optimisation of wastewater treatment plant control
988 to reduce greenhouse gas emissions. *Water Research* 55, 52–62.
- 989 Szemis, J.M., Maier, H.R., Dandy, G.C., 2012. A framework for using ant colony optimization to schedule
990 environmental flow management alternatives for rivers, wetlands, and floodplains. *Water Resources*
991 *Research* 48, 1–21.
- 992 Termansen, M., McClean, C.J., Preston, C.D., 2006. The use of genetic algorithms and Bayesian classifi-
993 cation to model species distributions. *Ecological Modelling* 192, 410–424.
- 994 Thuiller, W., 2003. BIOMOD - Optimizing predictions of species distributions and projecting potential
995 future shifts under global change. *Global Change Biology* 9, 1353–1362.
- 996 Thuiller, W., Lafourcade, B., Engler, R., Araújo, M.B., 2009. BIOMOD - A platform for ensemble fore-
997 casting of species distributions. *Ecography* 32, 369–373.
- 998 Tirelli, T., Pessani, D., 2009. Use of decision tree and artificial neural network approaches to model pres-
999 ence/absence of *Telestes muticellus* in piedmont (North-Western Italy). *River research and applications*
1000 25, 1001–1012.
- 1001 Van Broekhoven, E., Adriaenssens, V., De Baets, B., 2007. Interpretability-preserving genetic optimization
1002 of linguistic terms in fuzzy models for fuzzy ordered classification: An ecological case study. *International*
1003 *Journal of Approximate Reasoning* 44, 65–90.
- 1004 Van Broekhoven, E., Adriaenssens, V., De Baets, B., Verdonschot, P.F., 2006. Fuzzy rule-based macroin-
1005 vertebrate habitat suitability models for running waters. *Ecological Modelling* 198, 71–84.
- 1006 Vayghan, A.H., Zarkami, R., Sadeghi, R., Fazli, H., 2016. Modeling habitat preferences of Caspian kutum,
1007 *Rutilus frisii kutum* (Kamensky, 1901) (Actinopterygii, Cypriniformes) in the Caspian Sea. *Hydrobiologia*
1008 766, 103–119.
- 1009 Verberk, W., Verdonschot, P., van Haaren, T., van Maanen, B., 2012. Milieu-en habitatpreferenties van
1010 Nederlandse zoetwatermacrofauna. Technical Report. STOWA. Eindhoven.
- 1011 Verbyla, D.L., 1987. Classification trees: a new discrimination tool. *Canadian Journal of Forest Research*
1012 17, 1150– 1152.

- 1013 Vezza, P., Muñoz-Mas, R., Martínez-Capel, F., Mouton, A.M., 2015. Random forests to evaluate biotic
1014 interactions in fish distribution models. *Environmental Modelling & Software* 67, 173–183.
- 1015 VMM, 2018. Flemish Environment Agency. URL: <https://www.vmm.be>. Accessed on 2018-09-20.
- 1016 Weise, T., 2009. *Global Optimization Algorithms: Theory and Application*. volume 1.
- 1017 Whigham, P.A., 2000. Induction of a marsupial density model using genetic programming and spatial
1018 relationships. *Ecological Modelling* 131, 299–317.
- 1019 Wolpert, D.H., Macready, W.G., 1997. No free lunch theorems for optimization. *IEEE Transactions on*
1020 *Evolutionary Computation* 1, 67–82.
- 1021 Yee, T.W., Mitchell, N.D., 1991. Generalized additive models in plant ecology. *Journal of Vegetation*
1022 *Science* 2, 587–602.
- 1023 Zarkami, R., Sadeghi, R., Goethals, P.L.M., 2012. Use of fish distribution modelling for river management.
1024 *Ecological Modelling* 230, 44–49.
- 1025 Zarkami, R., Sadeghi, R., Goethals, P.L.M., 2014. Modelling occurrence of roach “*Rutilus rutilus*” in
1026 streams. *Aquatic Ecology* 48, 161–177.
- 1027 Zecchin, A.C., Simpson, A.R., Maier, H.R., Leonard, M., Roberts, A.J., Berrisford, M.J., 2006. Application
1028 of two ant colony optimisation algorithms to water distribution system optimisation. *Mathematical and*
1029 *Computer Modelling* 44, 451–468.
- 1030 Zuur, A.F., Ieno, E.N., Elphick, C.S., 2010. A protocol for data exploration to avoid common statistical
1031 problems. *Methods in Ecology and Evolution* 1, 3–14.
- 1032 Zuur, A.F., Ieno, E.N., Walker, N.J., Saveliev, A.A., Smith, G.M., 2009. *Mixed Effects Models and*
1033 *Extensions in Ecology with R*. 1 ed., Springer Science + Business Media, New York.

Box 1: Terminology machine learning

Machine learning: Field of research using computer programs or algorithms that have the ability to adapt or change (*i.e.* train) from an experience (data) given a performance measure.

Supervised learning: Learning from an experience by providing the machine learning algorithm with a set of inputs together with the corresponding outputs (labels).

Unsupervised learning: Learning from an experience by providing the machine learning algorithm with only inputs (and no labels) or in other words, making the algorithm search for patterns in data without having any labels to test it.

Training data: Observations of an experience used by the machine learning algorithm to train a model.

Testing data: Observations of an experience used to test the trained model. It is important to note test data are not used during the training phase.

Feature space: Set of all possible combinations of features of the input dataset. Both training and testing data are samples from the input dataset.

Feature selection or (input) variable selection: Process of selecting a subset of *relevant* features. The objectives of feature selection are to avoid the risk of overfitting by reducing model complexity (1), improve cluster detection (2) and reduce computational cost (3).

Ensemble learning: Training with multiple machine learning methods to obtain better predictive performance than by only using one machine learning method.

Model parameters: Model elements that are internal to a model, whose value can be estimated with data and are context-dependent. The action of estimating parameters is referred to as ‘parameter estimation’, in which a unique set of model parameter values are estimated with data.

Hyperparameters: Machine learning algorithm free options that need to be set beforehand, determining the training strategy and related efficiency of the algorithm. The settings of these hyperparameters can influence an algorithm’s and a model’s performance. In this review, hyperparameters refer specifically to parameters related to the algorithm free options, whereas ‘parameters’ refer to model elements. To clarify the difference between parameters and hyperparameters, a number of examples are given in section 3.2. For a good mathematical introduction to hyperparameters and their role in machine learning, we refer to Bergstra and Bengio (2012).

Evolutionary algorithms (EAs) (or evolutionary computing): Metaheuristic search algorithms (strategy) inspired by processes observed in evolution, *i.e.* selection, crossover and mutation.

Table 1: Machine learning approaches developed and used in species distribution modelling. In the table, WOS = web of science, and C.n. = cumulative number of publications. A hyphen in the first column indicates that no acronyms/full names are found in the literature.

Approach / technique (acronym)	Short description	C.n. in WOS (08/11/2017)	Notable references
Artificial neural networks (ANNs)	Non-linear mapping structures inspired on the biological system of the brain.	2000: 6; 2010: 77; 2017: 166	Fukuda et al. (2013); D'heygere et al. (2006)
BIOCLIM (-)	Delineates a rectangular environmental (bioclimate) hyperspace (or envelope) to estimate the response of species to a number of bioclimatic input variable.	1990/2000: 2; 2010: 1; 2017: 33	Carpenter et al. (1993); Elith et al. (2006)
Biodiversity modelling (BIOMOD)	Ensemble modelling platform/software combining several techniques.	2010: 11; 2017: 53	Thuiller (2003); Thuiller et al. (2009)
CLIMEX (-)	Model based on GROWEST, using a growth and stress index	sporadically used before 1990	Sutherst and Maywald (1985)
Decision trees (DT)	Classifiers expressed as a recursive partition or tree of the feature space.	2000: 5; 2010: 126; 2017: 421	Iverson and Prasad (1998)
Fuzzy logic (FL)	Method that deals with linguistic uncertainty by generalizing classical logic.	2000: 2; 2010: 20; 2017: 61	Adriaenssens (2004); Van Broekhoven et al. (2006)
Genetic algorithm for rule set production (GARP)	EA-inspired method used to produce a rule-bank SDM.	2010: 12; 2017: 224	Peterson et al. (2002); Stockwell and Noble (1992)
Generalized linear models (GLMs)	Collection of parametric techniques based on a random component, a systematic component, and a link function describing a relation between the former the random and systematic component.	2000: 23; 2010: 230; 2017: 600	Nelder and Wedderburn (1972); Zuur et al. (2010)
Generalized additive models (GAMs)	Extension of GLMs which relate the response variable to a linear combination of smoother functions.	2000: 43; 2010: 313; 2017: 738	Zuur et al. (2010)
GROWEST (-)	Model using a growth index based temperature, light, moisture	sporadically used before 1980	Nix et al. (1977)
Maximum entropy method (Maxent)	Technique using the principle of maximum entropy to make predictions from incomplete knowledge.	2010: 145; 2017: 1391	Phillips et al. (2006); Elith et al. (2011)
Random forest (RF)	Technique using bootstrap aggregation to create a set of decision trees.	2010: 27; 2017: 283	Prasad et al. (2006); Cutler et al. (2007)

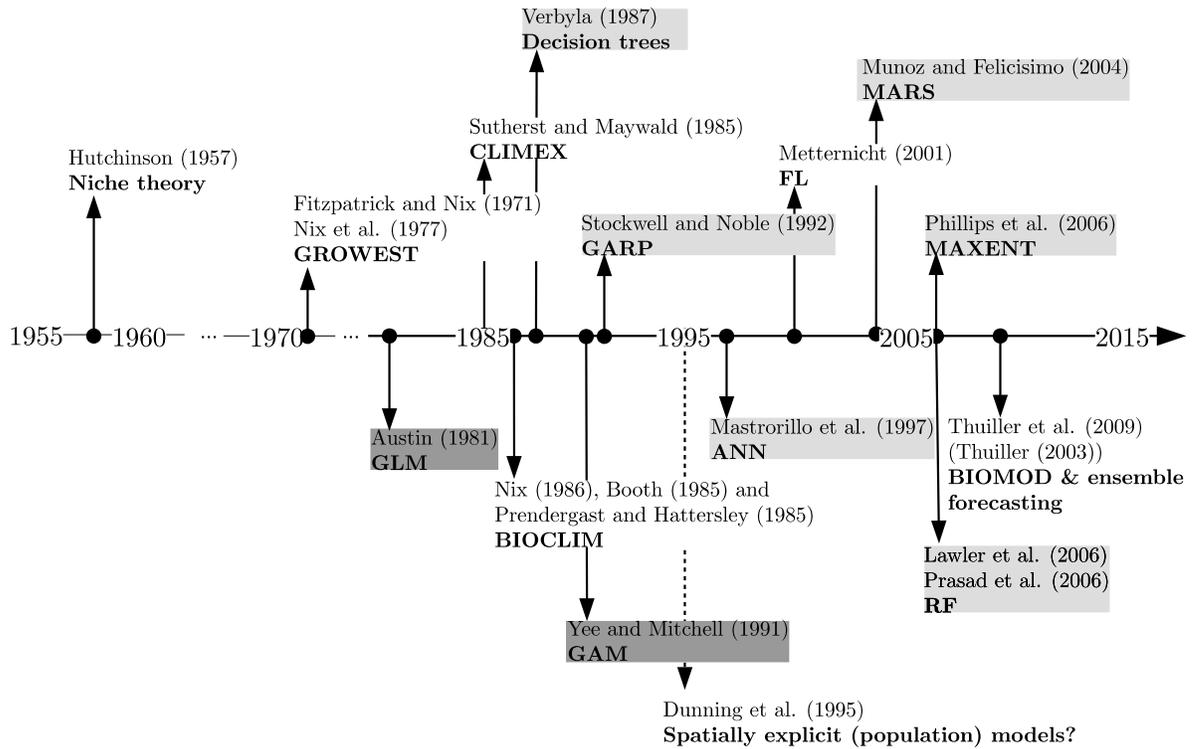


Figure 1: Chronological overview of first use or report of SDM modelling technique. Methods classified before 2000 under machine learning are indicated in a light grey boxes. Other methods that are often categorized under machine learning are shown in a dark grey boxes. This graph is a result of a web of science search on 08/11/2017 based on an intersection of literature on the shown methods and the field of ‘species distribution modelling’ (see supportive information 2). The acronyms found in this figure are listed in Table 1.

Box 2: Terminology EAs

Phenotype: Candidate solution to a problem, here represented by a model.

Genotype: Representation of a phenotype in a data type. Typically used genotypes are binary, real-valued strings or tree structures.

Genome: A specific formulation of the genotype (*e.g.* 1111011010).

Allele: A single element of the genotype (*e.g.* one bit).

Fitness: A measure of how good a solution to a problem is.

Chromosome: Object containing a genome and fitness.

Mapper: User-defined function which translates the genotype to phenotype.

Selection: Process of selecting chromosomes as parents for crossover, typically based on their fitness values.

Crossover: Process of combining the parents' genome to form genomes for the offspring.

Mutation: Process of randomly altering the parts of the genome.

Genetic algorithm: Evolutionary algorithms that use selection, crossover and mutation operators to solve an optimization problem. An explicit difference between GAs and other EAs is that GAs are designed as problem-independent solvers, whereas other EAs are designed to solve specific problems.

Table 2: Overview of literature review. The ‘subject of training’ is either parameter estimation, feature selection and/or hyperparameter optimization. Parameter estimation refers to the estimation of a unique set of model parameter values (see also Box 1). Hyperparameter optimization refers to the search for values of algorithm settings which influence an algorithm’s performance. ‘Training robustness’ indicates the robustness of the algorithm towards different samples of the data (by cross-validation, or bootstrapping) whereas ‘algorithm robustness’ is the robustness of the algorithm tested on the same data sample. The acronyms found in the column ‘objective function’ are found in Table 3. * = or no improvement 50 generations. ‘?’ = information was unclear or uncertain. Prev. = prevalence, *i.e.* number of species occurrence over number of samples. x = number of input features.

author; ecosystem	prev.	subject of training	objective functions	type of EA & operators	problem size	hyperparameters	hyperparameter optimization	training robustness	algorithm robustness	resampling scheme
D’Angelo et al. (1995); freshwater	?	parameter estimation; feature selection	ρ ; SSE	genetic algorithms genetic programming; mixed string genotype	$\pm 10^{10}$	# chromosomes: 200; crossover rate: 0.8; mutation rate: 0.1; # generations: 20000	iterative	no	10 runs	no
Whigham (2000); terrestrial	?	feature selection; parameter estimation	?	genetic programming; tree-genotype	?	?	?	?	100 runs	?
McKay (2001); terrestrial	?	feature selection; parameter estimation	?	genetic programming; tournament selection; half ramped initialization; tree-genotype	$\pm 10^{10}$	# chromosomes: 50; crossover rate: 0.9; mutation rate: 0.1; # generations: 50	no	?	?	?
D’hegyere et al. (2003); freshwater	?	feature selection	CCI	simple genetic algorithm; roulette wheel selection; binary string	2^{15}	# generations: 20; crossover rate: 0.6; mutation rate: 0.033; # generations: 20	?	10-fold cross-validation	?	?
McClellan et al. (2005); terrestrial	?	parameter estimation	AUC	Bayesian genetic algorithm; continuous string	$\pm 10^{2*9}$?	?	?	?	?
D’hegyere et al. (2006); freshwater	?	feature selection	CCI	simple genetic algorithm; roulette wheel selection; binary string	2^{17}	# chromosomes: 20; crossover rate: 0.6; mutation rate: 0.03; # generations: 40	iterative	10-fold cross-validation	x runs	stratified
Termansen et al. (2006); terrestrial	?	parameter estimation	AUC	Bayesian genetic algorithm; continuous string genotype	10^{2*9}	# chromosomes: 100; crossover rate: ?; mutation rate: <0.01; # generations: 90	iterative	?	x runs	?
Van Broekhoven et al. (2007); freshwater		parameter estimation	% CFCI	simple genetic algorithm; tournament selection; elitism; binary string and continuous	?	# chromosomes: 100; crossover rate: 0.95; mutation rate: $\sim 1/(\text{length chromosome})$ (<0.01); # generations: 1000*	iterative	?	100 runs	?
Fukuda and Hiramatsu (2008); freshwater	?	parameter estimation	MSE	simple genetic algorithm; binary string	?	?	?	10 runs	?	
Fukuda (2009); freshwater	?	parameter estimation	MSE	simple genetic algorithm; elitism	?	?	?	50 runs	?	
Tirelli and Pessani (2009); freshwater	0.7	feature selection	Kappa (?)	?	?	?	?		?	
Hoang et al. (2010); freshwater	0.12 to 0.72	feature selection	CCI; Kappa	binary string	2^{21}	?	?	3-fold cross-validation	5 runs	no
Favaro et al. (2011); freshwater	0.56	feature selection	Sn; Sp; CCI; Kappa; AUC	?	?	?	?	10-fold cross-validation	?	?

continued ...

... Table 2 continued

author; ecosystem	prev.	subject of training	objective functions	type of EA operators	problem size	hyperparameters	hyperparameter tuning	training robustness	algorithm robustness	resampling scheme
Fukuda et al. (2011); freshwater	0.27	parameter estimation	MSE	simple genetic algorithm; roulette wheel selection; elitism; binary string	2^{4*35}	# chromosomes: 100; crossover rate: ?; mutation rate: 0.05; # generations: 5000	?	3-fold cross-validation	20 runs	stratified based on prevalence
Jeong et al. (2011); terrestrial (marine)	?	feature selection; parameter estimation	RMSE	genetic algorithm genetic programming; tree-genotype	$\pm 10^{10}$	# chromosomes: 200; crossover rate: 0.6-0.9; mutation rate: 0-0.3; # generations: 100	?	bootstrapping	?	80 % training data 20 % test data
Fukuda et al. (2012); freshwater	0.5	parameter estimation	MSE	binary string	2^{4*17}	# chromosomes: 100; crossover rate: ?; mutation rate: ?; # generations: 2000	?	5-fold cross-validation	20 runs	?
Zarkami et al. (2012); freshwater	?	feature selection	CCI; Kappa	simple genetic algorithm; binary string	?	# chromosomes: 20; crossover rate: 0.6; mutation rate: 0.033; # generations: 20	?	3-fold cross-validation	?	?
Boets et al. (2013); freshwater	0.38	feature selection	Kappa; AUC; CCI	simple genetic algorithm; tournament selection; binary string	2^{11}	# chromosomes: 20 crossover rate: 0.6; mutation rate: 0.033; # generations: 20	iterative	3-fold cross-validation	?	random
Boets et al. (2013); freshwater	0.29	feature selection	Kappa; AUC; CCI	simple genetic algorithm; tournament selection; binary string	2^{11}	# chromosomes: 20; crossover rate: 0.6; mutation rate: 0.033; # generations: 20	iterative	3-fold cross-validation	?	random
Sadeghia et al. (2013); freshwater	?	feature selection	?	simple genetic algorithm; binary string	$\pm 10^{33}$	# generations: 20; crossover rate: 0.6; crossover rate: 0.033; # generations: 20	iterative	3-fold cross-validation	x runs	?
Sadeghi et al. (2014); freshwater	?	feature selection	?	binary string	$\pm 10^{33}$?	?	4-fold cross-validation	5 runs	?
Zarkami et al. (2014); freshwater	0.5	feature selection	CCI; Kappa	simple genetic algorithm; binary string	2^{10}	# chromosomes: 20; crossover rate: ?; mutation rate: ?; generations: 20	?	3-fold cross-validation	?	random
Muñoz-Mas et al. (2016a); freshwater	0.37	feature selection; hyperparameter optimization (decision tree)	TSS	GA with a derivative quasi-Newton method; mixed string genotype	$\pm 10^{2x}$	# chromosomes: 500; crossover rate: 0.75; mutation rate: 0.75; # generations: 500	?	3 times 3-fold cross-validation	?	stratified based on prevalence
Muñoz-Mas et al. (2016b); freshwater	0.62 0.66	feature selection	TSS penalty complexity	GA with a derivative quasi-Newton method; mixed string genotype	?	# chromosomes: 1000; crossover rate: 0.75; mutation rate: 0.75; # generations: 1000	?	3 times 3-fold cross-validation	?	?
Vayghan et al. (2016); freshwater	0.57	feature selection	?	binary string	2^9	# generations: 20; crossover rate: 0.6; mutation rate: 0.033; # generations: 20	?	?	?	?
Gobeyn and Goethals (2017); freshwater	?	feature selection parameter estimation	AUC Kappa	simple genetic algorithm tournament selection variable length binary string	$> 2^{112}$	# chromosomes: 100; crossover rate: 1.0; mutation rate: 0.05; # generations: 50 to 2000	iterative and guidelines of Gibbs et al. (2008)	bootstrapping	?	?
Muñoz-Mas et al. (2017); freshwater	0.21 0.42	feature selection	TSS stimulated overprediction	GA with a derivative quasi-Newton method; binary string	?	# chromosomes: function of ensemble size; crossover rate: 0.6; mutation rate: 0.6; # generations: function of ensemble size	iterative?	?	?	stratified based on prevalence

Table 3: Overview of often-used measures to define an objective function. For a in-depth review and formulation of regression coefficients, we refer to Mouton et al. (2010).

Measure	Classification (C) or regression (R)	Acronym Symbol	Reference
Correctly classified instances	C	CCI	Mouton et al. (2010)
Cohen's Kappa	C	Kappa	Mouton et al. (2010)
Sensitivity	C	Sn	Mouton et al. (2010)
Sensitivity	C	Sp	Mouton et al. (2010)
True skill statistic	C	TSS	Mouton et al. (2010)
Area under the receiver operator characteristic curve	C	AUC	Mouton et al. (2010)
Correctly fuzzy classified instances	C	% CFCI	Van Broekhoven et al. (2006)
(Root) mean squared errors	R	(R)MSE	Based on species density, see Fukuda (2009) Based on suitability, see Fukuda et al. (2012)
Sum of squared errors	R	SSE	Based on population size: D'Angelo et al. (1995)
Linear correlation	R	ρ	Based on population size: D'Angelo et al. (1995)

Table 4: Suggested EA or metaheuristic algorithm useful in species distribution modelling. Algorithms followed by a '*' are suggested in this review, but have yet to be tested, or have only tested in a number of experiments within ecology.

Learning problem	suggested algorithm	example reference
Feature selection	binary SGA	D'heygere et al. (2006)
	ant colony optimization*	-
Parameter estimation	real-coded SGA	Van Broekhoven et al. (2007)
	particle swarm optimization*	-
	simulated annealing*	-
Parameter estimation and feature selection	GP and GAs	Jeong et al. (2011)
	variable length GAs*	Gobeyn and Goethals (2017)
Hyperparameter optimization of other machine learning algorithms	evolutionary optimization (problem-specific)	Muñoz-Mas et al. (2016b)
Multi-objective optimization	NSGA-II*	-
Problem characteristics identification	EDA*	-

Box 3: General guidelines for applying a metaheuristic machine learning algorithm

Model:

1. **Model formulation:** First the model scale, resolution, model inputs, states, parameters and boundary conditions relevant to the species and case study need to be considered. At this stage, it is inspected if the model can be simplified by making specific assumptions (*e.g.* only consider specific species interactions) and/or identifying correlating features. For the latter, this can be done a priori model fitting with filtering methods based on input data (using, for instance, the Spearman rank correlation, see Saeys et al. (2007) and Dormann et al. (2012)) or during model fitting (by computing, for instance, a mutual information criterion, see May et al. (2008)). It is important to note that the use of automated procedures to select features should not serve as a replacement for an expert-based selection (Araújo and Guisan, 2006).
2. **Objective:** Define specific objectives and criteria to which the model should comply. For instance, is the aim of to obtain models which estimate primarily species presence well or rather species absence? Other aspects involved can be related to model complexity (for example see (Phillips et al., 2006)). Many options are available to define a measure. It is important to note that these are sensitive to the sample prevalence (Mouton et al., 2009, 2010). In other words: the measure used for model training can be varied as a function of the sample prevalence. This can cause a bias in the obtained model.
3. **Subject of training:** Defining which specific elements of the models are perturbed to maximize or minimize an objective function. It can be aimed to estimate model parameters or/and reduce the number of input variables (and thus model elements) (Fukuda et al., 2011). When the goal is to decrease the number of model structural elements, both wrapper and embedded feature selection methods can be used. feature selection selects input features which are most relevant - given an objective - to explain patterns in the data. In embedded feature selection, parameters are estimated while performing feature selection. On the contrary, in wrapper feature selection, parameters are estimated for each feature subset (or prior feature selection) (Saeys et al., 2007). For the latter, another machine learning algorithm is typically run within the feature selection procedure.

Box 3 (continued): General guidelines for applying a metaheuristic machine learning algorithm

Algorithm:

1. **Type of algorithm:** Users are advised to consider EAs and other metaheuristics in order to train models to test complex ecological hypotheses. A very good example using a simulated annealing metaheuristic algorithm to inspect the effect of species interactions and stress tolerance on biodiversity is presented by Baert et al. (2016). For relatively simple questions mainly aiming to get a first insight into the problem, we advise using Maxent, GLMs and/or decision trees.
2. **Encoding:** A binary (string of zeros and ones) encoding can be considered for wrapper feature selection. Applied to EAs, real-valued encoding (string of continuous values) can be considered for parameter estimation and a list of list encoding for embedded feature selection. Haupt and Haupt (2004) provide background hints and tips for the implementation of a binary and real-values encoding in evolutionary algorithms. To implement embedded feature selection in EAs, a list of list approach can be used (Srikanth et al., 1995; Gobeyn and Goethals, 2017). In addition, boundary conditions need to be addressed in the encoding and functioning of the operators (*e.g.* implementing repair operators for genome in EAs).
3. **Operators:** Metaheuristic algorithms use a number of operators depicting efficiency of the algorithm. Many implementations are available and depend on the encoding of the solutions. For EAs, typically tournament selection, uniform crossover and uniform mutation operators are implemented (Haupt and Haupt, 2004). The implementation of multiobjective machine learning with EAs requires specific selection operators (Deb et al., 2002).
4. **Hyperparameters:** All machine learning methods have a number of hyperparameters which need to be set. For standard EAs, selection, crossover and mutation rate needs to be set together with the number of iterations. Guidelines are available for many algorithms (*e.g.* EA, Gibbs et al. (2008) or decision trees, Everaert et al. (2016)). It is important to note that setting hyperparameters is for many machine learning algorithms required to acquire satisfying results with the lowest computational effort (Gibbs et al., 2015).

Box 3 (continued): General guidelines for applying a metaheuristic machine learning algorithm

Implementation:

1. **Programming language:** A number of programming languages exist to implement algorithms. For data science and machine learning, high-level programming languages such as Python and R are the most popular ones. These languages offer an interface for intuitive high-level programming. In addition, support can easily be found at online platforms such as Stack Overflow (www.stackoverflow.com) to solve specific programming problems. An alternative is to use algorithms available under a GUI environment, *e.g.* WEKA which also facilitates a command-line interface and Java API. These are helpful for machine learning, however, the use of these techniques to train hypothetical-driven SDMs can be tedious. In addition, these GUI applications are often difficult to use for repeated analysis (*i.e.* for uncertainty analysis) on high-performance computing infrastructure.
2. **Open science:** Developing a tailor-made package for a specific application can be a time-consuming practice. Therefore we recommend the development of an application based on existing Python or R packages. General or specific packages can be downloaded from the language developers websites and Github (<https://github.com/>). The latter is an open source code hosting platform for version control and (scientific) collaboration. Programmers in environmental and ecological science are increasingly aware of the importance of open source (science), code collaboration, reproducibility (not only in results but also in code) and modular scripting. A good example of this philosophy applied to species distribution modelling is published by Golding et al. (2017). In this approach, the authors provide a modular framework operating on snippets of R code that are interchangeable among each other. Finally, benchmarking algorithms and codes can be done by using datasets from GBIF, a free and open access to biodiversity data. In this case, an ecological data set is used, and different algorithms are applied to entangle the specific strengths and weaknesses of the used algorithms. As an alternative, an open, organized, online ecosystem for machine learning such as OpenML (<https://www.openml.org/home>) and Kaggle (<https://www.kaggle.com>) can be used.