1    **Automated discovery of relationships, models and principles in ecology**

2

3    **Running title:** Automated discovery in ecology

4

5    Word count: 6974

6

7    Pedro Cardoso[1,2,*], Paulo A. V. Borges[2], José C. Carvalho[2,3], François Rigal[2], Rosalina Gabriel[2],

8    José Cascalho[4,5], Luís Correia[5]

9

10    [1]*Finnish Museum of Natural History, University of Helsinki, P.O.Box 17 (Pohjoinen*

11    *Rautatiekatu 13), 00014 Helsinki, Finland.*

12    [2]*CE3C – Centre for Ecology, Evolution and Environmental Changes / Azorean Biodiversity*

13    *Group and Universidade dos Açores - Departamento de Ciências Agrárias, Rua Capitão João*

14    *d'Ávila, 9700-042 Angra do Heroísmo, Açores, Portugal. Email: pborges@uac.pt, rigal@uac.pt,*

15    *rgabriel@uac.pt.*

16    [3]*Department of Biology, CBMA – Molecular and Environmental Centre, University of Minho,*

17    *4710-057 Braga, Portugal. Email: josecarvalho@bio.uminho.pt.*

18    [4]*NIDes - Núcleo de Investigação e Desenvolvimento em e-Saúde, Departamento de Ciências*

19    *Agrárias, Rua Capitão João d'Ávila, 9700-042 Angra do Heroísmo, Açores, Portugal. Email:*

20    *jmc@uac.pt.*

21    [5]*BioISI – Biosystems and Integrative Sciences Institute, Faculdade de Ciências da Universidade*

22    *de Lisboa, Campo Grande, 1749-016 Lisboa, Portugal. Email: luis.correia@ciencias.ulisboa.pt.*

23    *\*Corresponding author: E-mail: pedro.cardoso@helsinki.fi.*

1

24

25

**ABSTRACT**

1. Ecological systems are the quintessential complex systems, involving numerous high-order interactions and non-linear relationships. The most commonly used statistical modelling techniques can hardly reflect the complexity of ecological patterns and processes. Finding hidden relationships in complex data is now possible through the use of massive computational power, particularly by means of Artificial Intelligence methods, such as evolutionary computation.

2. Here we use symbolic regression (SR), which searches for both the formal structure of equations and the fitting parameters simultaneously, hence providing the required flexibility to characterize complex ecological systems.

3. First, we demonstrate how SR can deal with complex datasets for: 1) modelling species richness; and 2) modelling species spatial distributions. Second, we illustrate how SR can be used to find general models in ecology, by using it to: 3) develop species richness estimators; and 4) develop the species-area relationship and the general dynamic model of oceanic island biogeography.

4. All the examples suggest that evolving free-form equations purely from data, often without prior human inference or hypotheses, may represent a very powerful tool for ecologists and biogeographers to become aware of hidden relationships and suggest general theoretical models and principles.

46

47    **Key-words:** artificial intelligence, evolutionary computation, genetic programming, species

48    richness estimation, species-area relationship, species distribution modelling, symbolic

49    regression.

50

51

52    **INTRODUCTION**

53

54    *Ecology as a complexity science*

55

56    Complexity is a term often used to characterize systems with numerous components interacting

57    in ways such that their collective behaviour is difficult to predict, but where emergent properties

58    give rise to, more or less simple but seldom linear, patterns (Table 1)(Holland 1995; Mitchell

59    2009). Complexity science is therefore an effort to understand non-linear systems with multiple

60    connected components and how "the whole is more than the sum of the parts" (Holland 1998).

61    Biological systems probably are among the most complex (Solé & Goodwin 2000), and among

62    them, ecological systems are the quintessential complex systems (Anand 2010). These are

63    composed of individuals, populations from different species, interacting and exchanging energy

64    in multiple ways, furthermore relating with the physical environment at different spatial and

65    temporal scales in non-linear relationships. As a consequence, ecology is dominated by

66    idiosyncratic results, with most ecological processes being contingent on the spatial and temporal

67    scales in which they operate, which makes it difficult to identify recurrent patterns, knowing also

68    that pattern does not necessarily identify process (Lawton 1996; Dodds 2009; Passy 2012). The

69    most commonly used exploratory (e.g. PCA, NMDS) and statistical modelling techniques (e.g.

70    linear and non-linear regression) can hardly reflect the complexity of ecological patterns and

71    processes, often failing to find meaningful relationships in data. More flexible techniques (e.g.

72    GAMs) usually do not allow an easy interpretation of results and particularly of putative causal

73    relationships. For ecological data, we require more flexible and robust, yet amenable to full interpretation,

74    analytical methods, which can eventually lead to the discovery of general principles and models.

75

76    *General principles and models in ecology*

77

78    The ultimate aim of any ecological principle is to provide a robust model for exploring,

79    describing and predicting ecological patterns and processes regardless of taxon identity and

80    geographic region (Lawton 1996; Dodds 2009). Finding a recurrently high goodness-of-fit for a

81    model to an ecological pattern for most taxa and ecosystems is usually the most compelling

82    evidence of a mechanistic process controlling that pattern. When general principles are translated

83    into robust models, general statistical methods are mostly abandoned in favor of these (Appendix

84    1). Such general, widely applicable, equations are mostly found by intellectual *tour de force*.

85    Yet, they surely are only the tip of the iceberg, usually incorporating few of the variables

86    increasingly available to ecologists and that could potentially explain such patterns.

87

88    *Computing power applied to complex ecological systems*

89

90    The automation of techniques for collecting and storing ecological and related data, with

91    increasing spatial and temporal resolutions, has become one of the central themes in ecology and

92    bioinformatics. Yet, automated and flexible ways to synthesise such complex and big data were

93    mostly lacking until recently. Finding hidden relations within such data is now possible through

4

94    the use of massive computational power. New computer-intensive methods have been developed

95    or are now available or possible (Reshef *et al.* 2011) including in particular the broad field of

96    Artificial Intelligence (AI) which has produced a variety of approaches. AI includes a series of

97    evolution-inspired techniques, brought together in the sub-field of evolutionary computation, of

98    which the most studied and well-known probably are genetic algorithms (Holland 1975). Genetic

99    programming, namely in the form of symbolic regression (SR)(Koza 1992), is a particular

100    derivation of genetic algorithms that searches the space of mathematical equations without any

101    constraints on their form, hence providing the required flexibility to represent complex systems

102    as presented by many ecological systems (Fig 1). Contrary to traditional statistical techniques,

103    symbolic regression searches for both the formal structure of equations and the fitting parameters

104    simultaneously (Schmidt & Lipson 2009). Finding the structure of equations is especially useful

105    to discover general models, providing general insights into the processes and eventually leading

106    to the discovery of new and as yet undiscovered principles. Fitting the parameters provides

107    insight into the specific data, and allow specific predictions.

108    Successful examples on the use of SR in ecology include modelling of land-use change (Manson

109    2005; Manson & Evans 2007), effects of climate change on populations (Tung *et al.* 2009;

110    Larsen *et al.* 2014), community distribution (Larsen, Field & Gilbert 2012; Yao *et al.* 2014),

111    predicting micro-organismal blooms (Muttil & Lee 2005; Jagupilla *et al.* 2015), deriving

112    vegetation indices (Almeida *et al.* 2015) and using parasites as biological tags (Barrett,

113    Kostadinova & Raga 2005).

114    In this work we explain, test and demonstrate the usefulness of SR in uncovering hidden

115    relationships within typical ecological datasets. First, we demonstrate how SR can deal with

116    complex datasets, namely for: 1) modelling species richness; and 2) modelling species spatial

117    distributions. Second, we illustrate how SR can be used to find general models in ecology, by

118    using it to: 3) develop species richness estimators; and 4) develop the species-area relationship

119    (SAR) and the general dynamic model of oceanic island biogeography (GDM).

120

121

122    **MATERIALS AND METHODS**

123

124    *Symbolic regression*

125

126    Symbolic regression works as a computational parallel to the evolution of species (Fig 1). A

127    population of initial equations is generated randomly by combining different building blocks,

128    such as the variables of interest (independent explanatory variables), algebraic operators (+, −, ÷,

129    ×), analytic function types (exponential, log, power, etc.), constants and other ways to combine

130    the data (e.g. Boolean or decision operators). Being random, these initial equations almost

131    invariably fail, but some are slightly better than others. All are then combined through crossover,

132    giving rise to new equations with characteristics from both parents. Equations with better fitness

133    (e.g. higher $r^2$) have a higher probability of recombining. To avoid new equations being bounded

134    by initially selected building blocks or quickly losing variability along the evolutionary process,

135    a mutation step (acting on any building block) is added to the process after crossover. After

136    multiple generations, an acceptable level of accuracy by some of the equations is often attained

137    and the researcher stops the process.

138    For this work we used the software Eureqa (Schmidt 2015). For each run, the software outputs a

139    list of equations along an error/complexity Pareto front, with the most accurate equation for each

140    level of complexity being shown (Fig 2). The Pareto front often presents an "elbow", where

141    near-minimum error meets near-minimum complexity. The equation in this inflection is closer to

142    the origin of both axes and is a good starting point for further investigation – if both axes are in

143    comparable qualitative scales. Often, however, this inflection point is not obvious and a single

144    formula is not clearly best. In such cases, weights can be given to each of them through Bayesian

145    statistics,  using indices that positively weight accuracy and negatively weight complexity, such

146    as Akaike's Information Criterion - AIC (Akaike 1974). However, in all cases it is important to

147    check all formulas along the Pareto front. Often equations or models that make immediate sense

148    to the specific question may not be detected by these automated methods.

149

150    ***Case-studies***

151

152    *Modelling species richness*

153    Modelling and mapping the species richness of high diversity taxa at regional to large scales is

154    often impossible without extrapolation from sampled to non-sampled sites. Here, we used an

155    endemic arthropod dataset collected in Terceira Island, Azores. Fifty-two sites were sampled

156    using pitfall traps for epigean arthropods (Cardoso *et al.* 2009), 13 in each of four land-use types:

157    natural forest, exotic forest, semi-natural pasture and intensively managed pasture. This dataset

158    was randomly divided into training and test data (26 sites each). We explained and predicted

159    species richness per site using elevation, slope, annual average temperature, annual precipitation

160    and an index of disturbance (Cardoso *et al.* 2013).

161    As the response variable was count data, Generalized Linear Models (GLM) and Generalized

162    Additive Models (GAM) with a Poisson error structure with log link were used. We used the

163    package MuMIn (Barton 2015) and the R environment (Team 2015) for multi-model inference

164    based on AICc values, using all variables plus all possible interactions for GLM. For GAM we

165    used package gam (Hastie 2015). For the SR search we used only algebraic and analytic

166    operators (+, −, ÷, ×, log, power), in this and all examples below, so that outputs could be most

167    easily interpreted. The $r^2$ goodness of fit was used as the fitness measure. As there was no clearly

168    best formula, AICc was used to choose a single equation along the Pareto front (Appendix 2).

169    Both $r^2$ and AICc were used to compare GLM and GAM with SR on the test dataset. Here and in

170    subsequent analyses, all models with a ΔAICc value < 2 (the difference between each model's

171    AICc and the lowest AICc) were considered as receiving equal statistical support.

172

173    *Modelling species distributions*

174    Species distribution modelling (SDM) is widely used to fill gaps in our knowledge on individual

175    species distributions. One of the general statistical methods used for SDM is logistic regression.

176    Among the multiple alternatives, the principle of maximum entropy (Maxent)(Phillips, Anderson

177    & Schapire 2006) has been found to be particularly robust (Elith *et al.* 2006).

178    We modelled the potential distribution of two endemic Azorean species in Terceira Island: the

179    rare forest click-beetle *Alestrus dolosus* (Coleoptera, Elateridae) and the abundant but mostly

180    forest restricted spider *Canariphantes acoreensis* (Araneae, Linyphiidae). Given the intrinsic

181    differences between methods, we had to use different background datasets. Maxent used the

182    environmental maps of the islands with a resolution of 100 m, from where it extracted pseudo-

183    absences. We then converted the probabilistic potential distribution maps to presence/absence

184    using the maximum value of training sensitivity plus specificity as the threshold as

185    recommended by Liu et al. (Liu *et al.* 2005). Logistic regression and SR used presence/absence

8

186  data from the 52 sampled sites. We used the package MuMIn (Barton 2015) and the R

187  environment (Team 2015) for multi-model inference of logistic regression based on AICc values.

188  In the SR run a step function was included, so that positive and negative values were converted

189  to presence and absence, respectively. Absolute error, reflecting the number of incorrect

190  classifications, was used as the fitness measure. As inflection points of the Pareto fronts were

191  clear, the best SR formula for each species was chosen based on them (Appendix 2). In all cases

192  only the training data (26 sites) were used for running the models. Logistic regression, Maxent

193  and SR were compared in their performance for predicting presence and absence of species on

194  the 26 test sites using the True Skill Statistic - TSS (Allouche, Tsoar & Kadmon 2006).

195

196  *Developing species richness estimators*

197  Several asymptotic functions have been used to estimate species richness (Soberon & Llorente

198  1993), including the Clench function (Clench 1979), the negative exponential function and the

199  rational function (Ratkowsky 1990) (Appendix 1). Our objective was to rediscover or eventually

200  find asymptotic models that would outperform them. Two independent datasets were used

201  resulting from exhaustive sampling for spiders in 1ha plots, performed by 8 collectors during 320

202  hours of sampling in a single hectare using five different methods. The training dataset was from

203  a mixed forest in Gerês (northern Portugal) and the test dataset was from a *Quercus* forest in

204  Arrábida (southern Portugal) (Cardoso *et al.* 2008a; Cardoso *et al.* 2008b).

205  Randomized accumulation curves for both sites were produced using the R package BAT

206  (Cardoso, Rigal & Carvalho 2015) (the package also includes both datasets). The true diversity

207  of each site was calculated as the average between different non-parametric estimators (Chao 1

208  and 2, Jackknife 1 and 2). Because the sampled diversity in the training dataset reached a very

9

209    high completeness but we wanted to simulate typically very incomplete sampling, datasets with

210    10, 20, 40, 80 and 160 randomly chosen samples were extracted and used, in addition to the

211    complete 320 samples dataset, as independent runs in SR. Squared error was used as the fitness

212    measure. Additionally, we imposed a strong penalty to non-asymptotic functions, although these

213    were still allowed in the search process. The weighted and non-weighted scaled mean squared

214    errors implemented in BAT (Cardoso, Rigal & Carvalho 2015) were used as accuracy measures.

215

216    *Developing the species-area relationship (SAR) and the general dynamic model of oceanic*

217    *island biogeography (GDM)*

218    One of the most studied examples of SARs is their application to island biogeography (ISAR).

219    The shape of ISARs has been modelled by many functions, but three of the simplest seem to be

220    preferred in most cases, the power, exponential and linear models (Triantis, Guilhaumon &

221    Whittaker 2012)(Appendix 1).

222    The general dynamic model of oceanic island biogeography was proposed to account for

223    diversity patterns within and across oceanic archipelagos as a function of area and age of the

224    islands (Whittaker, Triantis & Ladle 2008). Several different equations have been found to

225    describe the GDM, extending the different SAR models with the addition of a polynomial term

226    using island age and its square ($TT^2$), depicting the island's ontogeny. The first to be proposed

227    was an extension of the exponential model (Appendix 1)(Whittaker, Triantis & Ladle 2008), the

228    power model extensions following shortly after (Fattorini 2009; Steinbauer 2013).

229    Our objective was to test if we could re-discover and eventually refine existing models for the

230    ISAR and GDM from data alone. We used the Azores and Canary Islands spiders (Appendix

231    3)(Cardoso *et al.* 2010) as training data. To independently test the generality of models arising

10

232   from spider data, we used bryophyte data from the same archipelagos (Appendix 3)(Aranda *et al.*

233   2014). The area and maximum time since emergence of each island were used as explanatory

234   variables and the native species richness per island as the response variables. The $r^2$ value was

235   used as the fitness measure. The best SAR and GDM equations found by SR were chosen based

236   on the inspection of the Pareto front (Appendix 2), but looking also for interpretability of the

237   models. These were then compared with the existing models using AICc using the R package

238   BAT (Cardoso, Rigal & Carvalho 2015).

239

240

241   **RESULTS**

242

243   *Modelling species richness*

244

245   The model selected by GLM was:

246

247   $$S = e^{5.381 + 0.003432H - 0.001904P - 0.05257D}$$

248

249   ($r^2 = 0.744$, AICc = 30.793), where H = altitude, P = precipitation and D = disturbance. Yet, the

250   GLM model seems to be overfitting, as the results with the test data were considerably worse ($r^2$

251   = 0.146, AICc = 63.672). Overfitting also occurred with GAM, as the model was extremely good

252   for the training data ($r^2 = 0.930$, AICc = 8.643) yet much worse for testing data ($r^2 = -0.077$,

253   AICc = 85.601). The SR results performed worse than GLM or GAM with the training data, with

254   the formula chosen according to AICc being:

255

256 $$S = 0.673 + (8.696 - 0.002P)^{0.006H - 2.461}$$

257

258   ($r^2 = 0.641$, AICc = 43.050). However, the SR equation performs considerably better than GLM

259   or GAM with the test data ($r^2 = 0.289$, AICc = 62.354), revealing a higher generality of this

260   formula.

261

262   **Modelling species distributions**

263

264   The potential distribution models are relatively similar for *C. acoreensis* but show marked

265   differences for *A. dolosus* (Fig 3). Symbolic regression outperforms both other models for *A.*

266   *dolosus* and is as good as Maxent for *C. acoreensis*, with both outperforming LR (Table 2). The

267   SR models are not only the best, presenting maximum values for TSS, but are also the easiest to

268   interpret. *A. dolosus* is predicted to have adequate environmental conditions in all areas above

269   614m elevation, being restricted to pristine native forest. *C. acoreensis* can potentially be present

270   in all areas with disturbance values below 41.3, occurring not only in native forest but also in

271   adjacent semi-natural grassland and humid exotic forest. The LR and Maxent models used a

272   large number of explanatory variables for *A. dolosus*, yet performed worse on the test data than

273   did SR.

274

275   **Developing species richness estimators**

276

277   For the training dataset, one asymptotic model was found by SR (Appendix 2):

12

278

279
$$S = \frac{aQ}{b + Q}$$

280

281 where $a$ and $b$ are fitting parameters. This model is in fact the Clench model with a different

282 formulation (Appendix 1), where the asymptote is $a$. A second, slightly more complex but better

283 fitting, model was found for partial datasets with 40 or more samples:

284

285
$$S = \frac{c + aQ}{b + Q}$$

286

287 where $c$ is a third fitting parameter. The asymptote is again given by the value of $a$ (Fig 4). This

288 model is similar to the rational function (Appendix 1). It was found to outperform the Clench and

289 negative exponential for both the training and testing datasets (Table 3).

290

291 ***Developing the species-area relationship (SAR) and the general dynamic model of oceanic***

292 ***island biogeography (GDM).***

293

294 For the Azorean spiders, the best fitting previous model (both highest $r^2$ and lowest AICc) for the

295 ISAR was the exponential model (Table 4). The SR run discovered roughly the same model,

296 indicating, however, that the intercept ($c$ term) was adding unnecessary complexity. A similar

297 ranking of models was verified for bryophytes in the same region, revealing the robustness of the

298 new model.

299 For the Canary Islands, the best model for spiders was a linear function of area:

300

$$S = 75 + 0.047A$$

302

303     ($r^2 = 0.364$, AICc = 65.631). Although it is easy to interpret, the explained variance is relatively

304     low. The SR run reached a much higher explanatory power:

305

$$S = 112 + (-1.002^A)$$

307

308     ($r^2 = 0.806$, AICc = 57.320). In this case though, the equation is over-fitting to the few available

309     data (7 data points), as this function is erratic creating a biologically indefensible model. The

310     reason the ISAR is hard to model for the Canary Islands spiders is because we were missing the

311     major component Time (Cardoso *et al.* 2010). This is depicted by the GDM, of which the best of

312     the current equations was found to be the power model described by Fattorini (Fattorini

313     2009)(Table 4). Nevertheless, using SR we were able to find an improved, yet undescribed,

314     model (Table 4). This represents a general model expanding the linear SAR:

315

$$S = c + zA + xT - yT^2$$

317

318     When tested with Canarian bryophytes, this new formulation is almost as good as the power

319     model (Table 4).

320

321

322     **DISCUSSION**

323

324    Symbolic regression has the advantage over most standard regression methods (e.g. GLM) in

325    being fully flexible, allowing a much better fitting to data with similar interpretability to, for

326    example, a linear regression. SR also has one or more advantages over other, commonly used,

327    highly flexible regression (e.g. GAMs) or machine learning techniques (e.g. neural-networks):

328    (1) numerical, ordinal and categorical variables are easily combined; (2) redundant variables are

329    usually eliminated in the search process and only the most important are retained if anti-bloat

330    measures (intended to reduce the complexity of equations) are used; (3) the evolved equations

331    are human-readable and interpretable; and (4) solutions are easily applied to new data. Using SR

332    we were able to "distil" free-form equations and models that not only consistently outperform

333    but are more intelligible than the ones resulting from rigid methods such as GLM or "black-

334    boxes" such as Maxent. We were also able to re-discover and refine equations for estimating

335    species richness based on sampling curves and the ISAR and GDM from data alone.

336    All the examples presented in this work suggest that evolving free-form equations purely from

337    data, often without prior human inference or hypotheses, may represent a yet unexplored but

338    very powerful tool for ecologists and biogeographers, allowing the finding of hidden

339    relationships in data and suggesting new ideas to formulate general theoretical principles.

340

341    ***From particular relations to general principles***

342

343    Scientific fields such as physics rarely rely on general statistical inference methods such as linear

344    regression for hypothesis testing. The complexity of ecology made such methods an imperative

345    in most cases. The method now presented not only allows the discovery of relationships specific

346   to particular datasets, but also the finding of general models, globally applicable to multiple

347   systems of particular nature, as we tried to exemplify. As mentioned, SR is designed to optimize

348   both the form of the equations and the fitting parameters simultaneously. The fitting parameters

349   usually are specific to each dataset, but the form may give clues towards some general principle

350   (e.g. all archipelagos will follow an ISAR even if each archipelago will have its own c and z

351   values). Although this aspect has not been explored in this study, we suggest two ways of finding

352   general principles.

353   First, as was hinted by our estimators' example, one may independently analyse multiple datasets

354   from the same type of systems. From each dataset, one or multiple equations may arise. Many of

355   these will be similar in form even if the fitting parameters are different. Terms repeated in

356   several equations along the Pareto front or with different datasets tend to be meaningful (Schmidt

357   & Lipson 2009). We may then try to fit the most promising forms to all datasets optimizing the

358   fitting parameters to each dataset and look for which forms seem to have general value over all

359   data.

360   Second, one may simultaneously analyse multiple datasets from the same type of systems but

361   with a change to the general SR implementation. Instead of optimizing both form and fitting

362   parameters, the algorithm may focus on finding the best form, with fitting parameters being

363   optimized during the evaluation step of the evolution for each dataset independently. This

364   parameter optimization could be done with standard methods such as quasi-newton or simplex

365   (Nocedal & Wright 1999). To our knowledge, this approach has yet to be implemented, but it

366   would allow finding general models and possibly principles, independently of the idiosyncrasies

367   of each dataset.

368

16

369    *The need for human inference*

370

371    Our results show that an automated discovery system can identify meaningful relationships in

372    ecological data. Yet, as shown by our Canary Island spider SAR model, some equations might be

373    very accurate but overfit the data. As with any relationship finding, either automated or human,

374    correlation does not imply causation and spurious relationships are not only possible but

375    probable given complex enough data.

376    Although the method here presented is automated, it is part of a collaborative human–machine

377    effort. The possibility of exploiting artificial intelligence working together with human expertise

378    can be traced back to Engelbart (Engelbart 1962), where the term "augmented intelligence" was

379    coined to designate such collaboration. It has been subsequently developed and extended to

380    teamwork involving one or more artificial intelligence agents together with one or (many) more

381    humans, in diverse domains such as robotic teams (Yanco, Drury & Scholtz 2004) or collective

382    intelligence for evolutionary multi-objective optimization (Cinalli 2015).

383    In ecological problems, human knowledge may play a fundamental role: 1) in the beginning of

384    the process, when we must select input variables, building blocks and SR parameters; and 2) in

385    the interpretation and validation of equations. The choice of equations along a machine-

386    generated Pareto front should also take advantage of human expert knowledge to identify the

387    most interesting models to explain the data. The researcher might then decide to disregard,

388    accept or check equation validity using other methods.

389

390    *A priori knowledge*

391

17

392     To some extent, it is possible to select *a priori* the type of models the algorithm will search for

393     by selecting the appropriate variables and building blocks. Another way to take advantage of

394     previous knowledge is to use as part of the initial population of equations some, possibly

395     simpler, equations we know are related with the problem. For example, when searching for the

396     GDM we could have given the algorithm multiple forms of the ISAR to seed the search process.

397     This should be complemented with random equations to create the necessary variation for

398     evolution.

399

400     ***Fine-tuning the process***

401

402     The number of options in SR is immense. Population size is positively correlated with variability

403     of models and how well the search space is explored, but might considerably slow the search.

404     Mutation rates are also positively correlated with variability, but rates that are too high might

405     prevent the algorithm converging on the best models. The fitness measure depends on the

406     specific problem and the type of noise expected.

407     The number of generations to let the search run is entirely dependent on the problem complexity

408     and time available. Often the algorithm reaches some equation that makes immediate sense to the

409     researcher and the process can be immediately stopped for further analysis of results. Sometimes

410     several equations seem to make sense but are not entirely convincing, in which case several

411     indicators can be used as a stop rule, such as high values of stability and maturity of the

412     evolution process (Schmidt 2015).

413     The speed with which evolution occurs is extremely variable, depending on factors including the

414     complexity of the relationships, having the appropriate variables and building blocks and the

18

415    level of noise in the data. Fortunately, the process is easily adaptable to parallel computing, as

416    many candidate functions can be evaluated simultaneously, allowing the use of multiple cores

417    and even computer clusters to speed the search of equations.

418

419    ***Caveats and alternatives***

420

421    The SR approach is fully data-driven. This means it requires high-quality data if meaningful

422    relationships are to be found. Also, it makes no a priori assumptions, so the final result might

423    make no (obvious) sense, leading to spurious inferences, particularly if data are scarce or poor-

424    quality, or if the right building blocks are not provided. Additionally, SR suffers from the same

425    limitations of evolutionary algorithms in general. In many cases the algorithm may get stuck in

426    local minima of the search space, requiring time (or even a restart with different parameters) to

427    find the global minimum.

428    Many data mining techniques are regarded, and rightly so, as "black boxes". SR is transparent in

429    this regard, as variables are related through human-interpretable formulas. This is particularly

430    important if the goal is to find equations with both predictive and explanatory power, building

431    the bridge between finding the pattern and explaining the driving process, or if a general

432    principle is to be suggested.

433

434    ***The automation of science?***

435

436    The methods here presented can be powerful additions to theoretical and experimental ecology,

437    even if new conceptual hypotheses have to be created to accommodate the new equations. Such

19

438  models could even be the only available means of investigating complex ecological systems

439  when experiments are not feasible or datasets get too big/complex to model, using traditional

440  statistical techniques.

441  This kind of techniques has led several authors to suggest the "automation of science" (King *et*

442  *al.* 2009), where computers are able to advance hypotheses, test them and reach conclusions in

443  largely unassisted processes. The SR potential as an exploratory step, to be reasoned alongside

444  and proven with other methods is also exciting. The resulting formulas will help researchers to

445  focus on initially imperceptible but interesting relationships within datasets and help guide the

446  process of hypothesis creation.

447

448

449  **ACKNOWLEDGEMENTS**

450

456

457

458  **DATA ACCESSIBILITY**

459  All data available in Appendix 2.

460

20

461

**REFERENCES**

463

Akaike, H. (1974) New Look at Statistical-Model Identification. *Ieee Transactions on Automatic Control,* **Ac19,** 716-723.

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

Almeida, J., dos Santos, J.A., Miranda, W.O., Alberton, B., Morellato, L.P.C. & Torres, R.D. (2015) Deriving vegetation indices for phenology analysis using genetic programming. *Ecological Informatics,* **26,** 61-69.

Anand, M., Gonzalez, A., Guichard, F., Kolasa, J. & Parrott, L. (2010) Ecological systems as complex systems: challenges for an emerging science. *Diversity,* **2,** 395-410.

Aranda, S.C., Gabriel, R., Borges, P.A.V., Santos, A.M.C., de Azevedo, E.B., Patino, J., Hortal, J. & Lobo, J.M. (2014) Geographical, Temporal and Environmental Determinants of Bryophyte Species Richness in the Macaronesian Islands. *Plos One,* **9**.

Barrett, J., Kostadinova, A. & Raga, J.A. (2005) Mining parasite data using genetic programming. *Trends in Parasitology,* **21,** 207-209.

Barton, K. (2015) MuMIn: Multi-Model Inference.

Cardoso, P., Aranda, S.C., Lobo, J.M., Dinis, F., Gaspar, C. & Borges, P.A.V. (2009) A spatial scale assessment of habitat effects on arthropod communities of an oceanic island. *Acta Oecologica-International Journal of Ecology,* **35,** 590-597.

Cardoso, P., Arnedo, M.A., Triantis, K.A. & Borges, P.A.V. (2010) Drivers of diversity in Macaronesian spiders and the role of species extinctions. *Journal of Biogeography,* **37,** 1034-1046.

484     Cardoso, P., Gaspar, C., Pereira, L.C., Silva, I., Henriques, S.S., da Silva, R.R. & Sousa, P. (2008a)

485          Assessing spider species richness and composition in Mediterranean cork oak forests. *Acta*

486          *Oecologica-International Journal of Ecology,* **33,** 114-127.

487     Cardoso, P., Rigal, F. & Carvalho, J.C. (2015) BAT - Biodiversity Assessment Tools, an R package for the

488          measurement and estimation of alpha and beta taxon, phylogenetic and functional diversity.

489          *Methods in Ecology and Evolution,* **6,** 232-236.

490     Cardoso, P., Rigal, F., Fattorini, S., Terzopoulou, S. & Borges, P.A.V. (2013) Integrating Landscape

491          Disturbance and Indicator Species in Conservation Studies. *Plos One,* **8**.

492     Cardoso, P., Scharff, N., Gaspar, C., Henriques, S.S., Carvalho, R., Castro, P.H., Schmidt, J.B., Silva, I.,

493          Szuts, T., De Castro, A. & Crespo, L.C. (2008b) Rapid biodiversity assessment of spiders (Araneae)

494          using semi-quantitative sampling: a case study in a Mediterranean forest. *Insect Conservation and*

495          *Diversity,* **1,** 71-84.

496     Cinalli, D., Martí, L., Sanchez-Pi, N., & Garcia, A.C.B. (2015) Collective preferences in evolutionary multi-

497          objective optimization: techniques and potential contributions of collective intelligence. *30th*

498          *Annual ACM Symposium on Applied Computing*, pp. 133-138.

499     Clench, H. (1979) How to make regional lists of butterflies: some thoughts. *Journal of the Lepidopterists'*

500          *Society,* **33,** 216-231.

501     Dodds, W.K. (2009) *Laws, theories and patterns in ecology*. University of California Press, Berkeley, CA.

502     Elith, J., Graham, C.H., Anderson, R.P., Dudik, M., Ferrier, S., Guisan, A., Hijmans, R.J., Huettmann, F.,

503          Leathwick, J.R., Lehmann, A., Li, J., Lohmann, L.G., Loiselle, B.A., Manion, G., Moritz, C.,

504          Nakamura, M., Nakazawa, Y., Overton, J.M., Peterson, A.T., Phillips, S.J., Richardson, K.,

505          Scachetti-Pereira, R., Schapire, R.E., Soberon, J., Williams, S., Wisz, M.S. & Zimmermann, N.E.

506          (2006) Novel methods improve prediction of species' distributions from occurrence data.

507          *Ecography,* **29,** 129-151.

508     Engelbart, D. (1962) Augmenting human intellect: a conceptual framework. (ed. S.R. AFOSR-3233).

509          Stanford Research Institute, Menlo Park, CA.

22

510    Evans, M.R., Grimm, V., Johst, K., Knuuttila, T., de Langhe, R., Lessells, C.M., Merz, M., O'Malley, M.A.,

511        Orzack, S.H., Weisberg, M., Wilkinson, D.J., Wolkenhauer, O. & Benton, T.G. (2013) Do simple

512        models lead to generality in ecology? *Trends in Ecology & Evolution,* **28,** 578-583.

513    Fattorini, S. (2009) On the general dynamic model of oceanic island biogeography. *Journal of*

514        *Biogeography,* **36,** 1100-1110.

515    Hastie, T. (2015) gam: Generalized Additive Models. R package version 1.12.

516    Holland, J.H. (1975) *Adaptation in natural and artificial systems : an introductory analysis with*

517        *applications to biology, control, and artificial intelligence*. University of Michigan Press, Ann

518        Arbor.

519    Holland, J.H. (1995) *Hidden order : how adaptation builds complexity*. Addison-Wesley, Reading, Mass.

520    Holland, J.H. (1998) *Emergence : from chaos to order*. Addison-Wesley, Reading, Mass.

521    Jagupilla, S.C.K., Vaccari, D.A., Miskewitz, R., Su, T.L. & Hires, R.I. (2015) Symbolic Regression of

522        Upstream, Stormwater, and Tributary E-Coli Concentrations Using River Flows. *Water*

523        *Environment Research,* **87,** 26-34.

524    King, R.D., Rowland, J., Oliver, S.G., Young, M., Aubrey, W., Byrne, E., Liakata, M., Markham, M., Pir,

525        P., Soldatova, L.N., Sparkes, A., Whelan, K.E. & Clare, A. (2009) The Automation of Science.

526        *Science,* **324,** 85-89.

527    Koza, J.R. (1992) *Genetic programming : on the programming of computers by means of natural selection*.

528        MIT Press, Cambridge, Mass.

529    Larsen, P.E., Cseke, L.J., Miller, R.M. & Collart, F.R. (2014) Modeling forest ecosystem responses to

530        elevated carbon dioxide and ozone using artificial neural networks. *Journal of Theoretical Biology,*

531        **359,** 61-71.

532    Larsen, P.E., Field, D. & Gilbert, J.A. (2012) Predicting bacterial community assemblages using an artificial

533        neural network approach. *Nature Methods,* **9,** 621-+.

534    Lawton, J.H. (1996) Patterns in ecology. *Oikos,* **75,** 145-147.

23

535     Liu, C.R., Berry, P.M., Dawson, T.P. & Pearson, R.G. (2005) Selecting thresholds of occurrence in the

536         prediction of species distributions. *Ecography,* **28,** 385-393.

537     Manson, S.M. (2005) Agent-based modeling and genetic programming for modeling land change in the

538         Southern Yucatan Peninsular Region of Mexico. *Agriculture Ecosystems & Environment,* **111,** 47-

539         62.

540     Manson, S.M. & Evans, T. (2007) Agent-based modeling of deforestation in southern Yucatan, Mexico,

541         and reforestation in the Midwest United States. *Proceedings of the National Academy of Sciences*

542         *of the United States of America,* **104,** 20678-20683.

543     Mitchell, M. (2009) *Complexity : a guided tour*. Oxford University Press, Oxford England ; New York.

544     Muttil, N. & Lee, J.H.W. (2005) Genetic programming for analysis and real-time prediction of coastal algal

545         blooms. *Ecological Modelling,* **189,** 363-376.

546     Nocedal, J. & Wright, S.J. (1999) *Numerical optimization*. Springer, New York.

547     Passy, S.I. (2012) A hierarchical theory of macroecology. *Ecology Letters,* **15,** 923-934.

548     Phillips, S.J., Anderson, R.P. & Schapire, R.E. (2006) Maximum entropy modeling of species geographic

549         distributions. *Ecological Modelling,* **190,** 231-259.

550     Ratkowsky, D.A. (1990) *Handbook of nonlinear regression models*. M. Dekker, New York.

551     Reshef, D.N., Reshef, Y.A., Finucane, H.K., Grossman, S.R., McVean, G., Turnbaugh, P.J., Lander, E.S.,

552         Mitzenmacher, M. & Sabeti, P.C. (2011) Detecting Novel Associations in Large Data Sets. *Science,*

553         **334,** 1518-1524.

554     Russell, S.J., Norvig, P. & Davis, E. (2010) *Artificial intelligence : a modern approach,* 3rd edn. Prentice

555         Hall, Upper Saddle River.

556     Schmidt, M. & Lipson, H. (2009) Distilling Free-Form Natural Laws from Experimental Data. *Science,*

557         **324,** 81-85.

558     Schmidt, M.L., H. (2015) Eureqa.

559     Soberon, J. & Llorente, J. (1993) The Use of Species Accumulation Functions for the Prediction of Species

560         Richness. *Conservation Biology,* **7,** 480-488.

561    Solé, R.V. & Goodwin, B.C. (2000) *Signs of life : how complexity pervades biology*. Basic Books, New

562        York.

563    Steinbauer, M.J., Klara, D., Field, R., Reineking, B. & Beierkuhnlein, C. (2013) Re-evaluating the general

564        dynamic theory of oceanic island biogeography. *Frontiers of Biogeography,* **5,** 185-194.

565    Team, R.D.C. (2015) R: A Language and Environment for Statistical Computing. R Foundation for

566        Statistical Computing, Vienna, Austria.

567    Triantis, K.A., Guilhaumon, F. & Whittaker, R.J. (2012) The island species-area relationship: biology and

568        statistics. *Journal of Biogeography,* **39,** 215-231.

569    Tung, C.P., Lee, T.Y., Yang, Y.C.E. & Chen, Y.J. (2009) Application of genetic programming to project

570        climate change impacts on the population of Formosan Landlocked Salmon. *Environmental*

571        *Modelling & Software,* **24,** 1062-1072.

572    Whittaker, R.J., Triantis, K.A. & Ladle, R.J. (2008) A general dynamic theory of oceanic island

573        biogeography. *Journal of Biogeography,* **35,** 977-994.

574    Yanco, H.A., Drury, J.L. & Scholtz, J. (2004) Beyond usability evaluation: Analysis of human-robot

575        interaction at a major robotics competition. *Human-Computer Interaction,* **19,** 117-149.

576    Yao, M.J., Rui, J.P., Li, J.B., Dai, Y.M., Bai, Y.F., Hedenec, P., Wang, J.M., Zhang, S.H., Pei, K.Q., Liu,

577        C., Wang, Y.F., He, Z.L., Frouz, J. & Li, X.Z. (2014) Rate-specific responses of prokaryotic

578        diversity and structure to nitrogen deposition in the Leymus chinensis steppe. *Soil Biology &*

579        *Biochemistry,* **79,** 81-90.

580

581

582 **Table 1. Glossary of terms.**

| |
|---|
| **Artificial intelligence (AI)** - A scientific field concerned with the automation of activities we associate with human thinking (Russell, Norvig & Davis 2010). |
| **Big data** - Very large amount of structured or unstructured data, hard to model with general statistical techniques but with the potential to be mined for information. |
| **Complex system** - A system in which a large network of components organize, without any central controller and simple although non-linear rules of operation, into a complex collective behaviour that creates patterns, uses information, and, in some cases, evolves and learns (Mitchell 2009). |
| **General model** - An equation that is found to be useful for multiple datasets, often but not necessarily, derived from a general principle. In most cases the formal structure of equations is kept fixed, while some parameters must be fitted for each individual dataset. |
| **General principle** - Refers to concepts or phenomenological descriptions of processes and interactions (Evans *et al.* 2013). May not have direct translation to any general model, but be a purely conceptual abstraction. |
| **Genetic programming (GP)** - A biologically-inspired method for getting computers to automatically create a computer program to solve a given problem (Koza 1992). It is a type of evolutionary algorithm, where each solution to be tested (individual in a population of possible solutions) is a computer program. |
| **Pareto front** - A curve connecting a set of best solutions in a multi-objective optimization problem. If several conflicting objectives are sought (e.g. minimize both error and complexity of formulas), the Pareto front allows visualizing the set of best solutions. |

> **Symbolic regression (SR)** - A function discovery approach for modelling of multivariate data. It is a special case of genetic programming, one where possible solutions are equations instead of computer programs.

583

584

585 **Table 2. Species distribution models for two endemic arthropod species on the island of**

586 **Terceira (Azores, Portugal).**

| Model | Formula | Sensitivity | Specificity | TSS |
|---|---|---|---|---|
| *Alestrus dolosus* | | | | |
| Logistic regression | $$\frac{1}{1+e^{-(8469-0.432P-540.7T)}}$$ | 0 | **1** | 0 |
| Maxent | Uses all variables but *Sl*, main is *D* (contribution = 74.1%) | 0.5 | **1** | 0.5 |
| Symbolic regression | $step(H - 614)$ | **1** | 0.75 | **0.75** |
| *Canariphantes acoreensis* | | | | |
| Logistic regression | $$\frac{1}{1+e^{-(3.617-0.103D)}}$$ | 0.667 | **0.7** | 0.367 |
| Maxent | Uses only *D* (contribution = 100%) | **0.833** | 0.65 | **0.483** |
| Symbolic regression | $step(41.3 - D)$ | **0.833** | 0.65 | **0.483** |

587 Accuracy statistics on an independent test dataset are given by the True Skill Statistic (TSS). $H =$

588 altitude, $Sl =$ slope, $T =$ average annual temperature, $P =$ annual precipitation and $D =$

589 disturbance index. The step function in symbolic regression converts positive values inside

590 parentheses to presence and negative values to absence. Best values in bold.

591

592

28

593 **Table 3. Comparison of three asymptotic equations used to estimate spider species richness**

594 **in two forest sites.**

| Model | Raw accuracy | Weighted accuracy |
|---|---|---|
| **Gerês (training)** | | |
| Observed | 0.113 | 0.037 |
| Clench | 0.055 | 0.018 |
| Negative exponential | 0.115 | 0.049 |
| Rational function | **0.045** | **0.012** |
| **Arrábida (testing)** | | |
| Observed | 0.103 | 0.031 |
| Clench | 0.038 | 0.010 |
| Negative exponential | 0.092 | 0.037 |
| Rational function | **0.032** | **0.008** |

595 See Appendix 1 for formulas. Raw accuracy is the scaled mean squared error considering the

596 entire observed accumulation curve (each formula was fitted to the curves using 4 to 320

597 samples) and weighted accuracy is this value weighted by the sampling effort at each point in the

598 curve (where effort is the ratio between number of individuals and observed species richness).

599 Note that lower values (in bold) are better as they reflect the deviation from a perfect estimator.

600

29

601   **Table 4. Species area relationship (SAR) models for Azorean taxa and General Dynamic**

602   **Models (GDM) of oceanic island biogeography for Canarian taxa.**

| Model | Formula | $r^2$ | AICc |
|---|---|---|---|
| **SAR Azorean Spiders (training)** | | | |
| Power | $S = 13.379 * A^{0.438}$ | 0.642 | 32.505 |
| Exponential | $S = 0.549 + 4.538 \log A$ | **0.780** | 28.102 |
| Linear | $S = 19.357 + 0.017A$ | 0.435 | 36.604 |
| Exponential (SR) | $S = 4.641 \log A$ | **0.780** | **23.319** |
| **SAR Azorean Bryophytes (testing)** | | | |
| Power | $S = 181.625 * A^{0.803}$ | 0.666 | 78.085 |
| Exponential | $S = - 27.824 + 57.114 \log A$ | **0.728** | 76.208 |
| Linear | $S = 196.215 + 0.259A$ | 0.617 | 79.295 |
| Exponential (SR) | $S = 51.889 \log A$ | 0.722 | **71.617** |
| **GDM Canarian Spiders (training)** | | | |
| Whittaker | $S = -185.589 + 41.732\log A + 17.776T -1.022T^2$ | 0.873 | 110.350 |
| Fattorini | $\log S = 2.585 + 0.281\log A + 0.157T -0.009T^2$ | 0.941 | 105.025 |
| Steinbauer | $\log S = 3.367 + 0.098\log A + 1.502\log T - 0.454\log T^2$ | 0.814 | 113.007 |
| SR | $S = 42.283 + 0.051A + 17.379T - T^2$ | **0.952** | **61.505** |
| **GDM Canarian Bryophytes (testing)** | | | |
| Whittaker | $S = -176.599 + 66.602\log A + 21.361T -1.620T^2$ | **0.773** | **125.214** |
| Fattorini | $\log S = 4.544 + 0.137\log A + 0.126T -0.009T^2$ | **0.803** | **124.217** |
| Steinbauer | $\log S = 5.136 + 0.017\log A + 1.063\log T - 0.382\log T^2$ | 0.612 | 128.963 |
| SR | $S = 192.660 + 0.075A + 20.702T - 1.576T^2$ | **0.785** | **124.841** |

603    S = native species richness, A = area of the island and T = maximum time of emergence. Best

604    models are indicated in bold.

605

606 **Fig. 1. Schematic representation of the symbolic regression workflow.**

607 The basic representation is a parse-tree where building blocks such as variables (in this case: $x_1$,

608 $x_2$), parameters (integers or real numbers) and operators (e.g. $+, -, \times, \div$) are connected forming

609 functions (in parenthesis under the first line of trees). Initial equations are generated by randomly

610 linking different building blocks. Equations are combined through crossover, giving rise to new

611 equations with characteristics from both parents (arrows linking the first and second rows of

612 trees). Equations with better fitness (e.g. $r^2$) have higher probabilities of recombining. To avoid

613 loss of variability, a mutation step is added after crossover (arrows linking the second and third

614 rows of trees). After multiple generations, evolution stops and a set of free-form equations best

615 reflecting the input data is found.

616

617 **Fig. 2. Example of a Pareto front depicting error vs. complexity.**
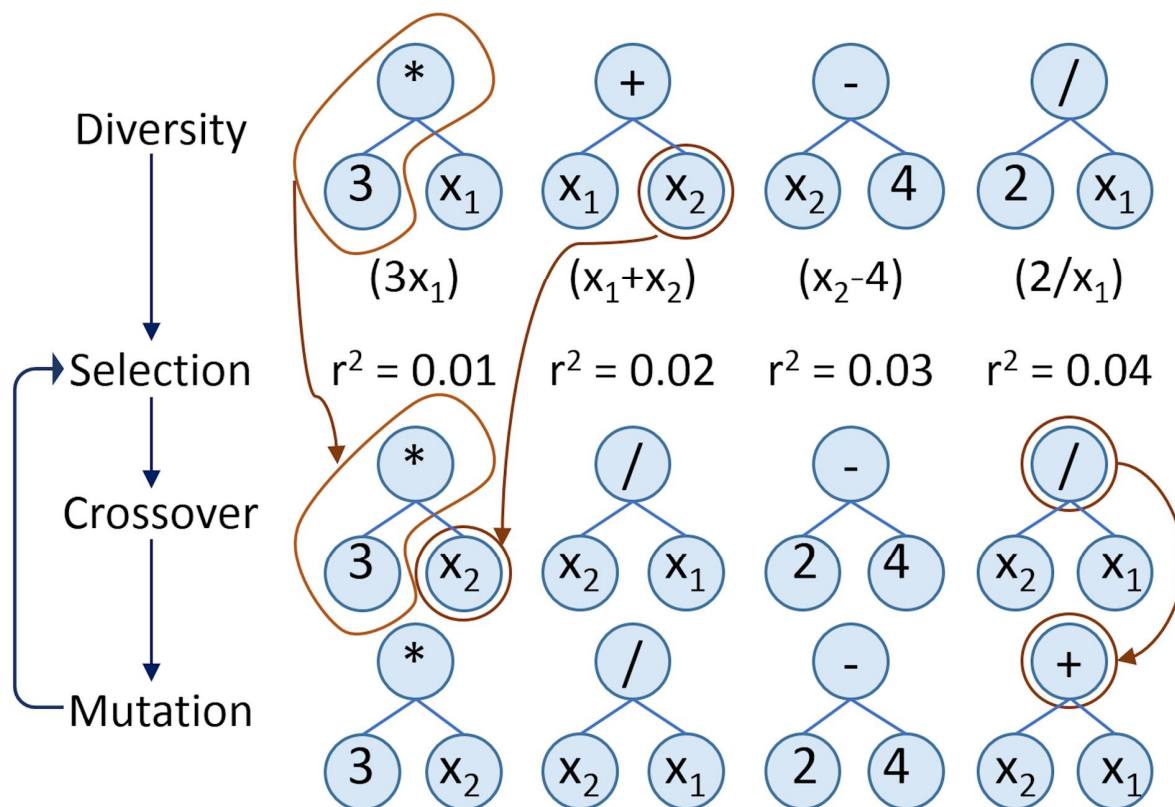
618 This example reflects a symbolic regression search of the best species–area relationship for native spiders

619 in the Azores (Portugal). The second formula is clearly the most promising, with both high accuracy (low

620 error) and low complexity. In many occasions a single formula is not clearly best, in which case

621 weights can be given to each of them through Bayesian statistics and multiple formulas presented

622 as possible outcomes.

623

624

625 **Fig. 3. Predicted distribution of two Azorean arthropods using three modelling methods.**

626 Observed locations (white dots) and predicted distribution (dark green areas) of *Alestrus dolosus*

627 (Coleoptera) and *Canariphantes acoreensis* (Araneae) in the island of Terceira (Azores,

628 Portugal) using logistic regression, maximum entropy and symbolic regression.
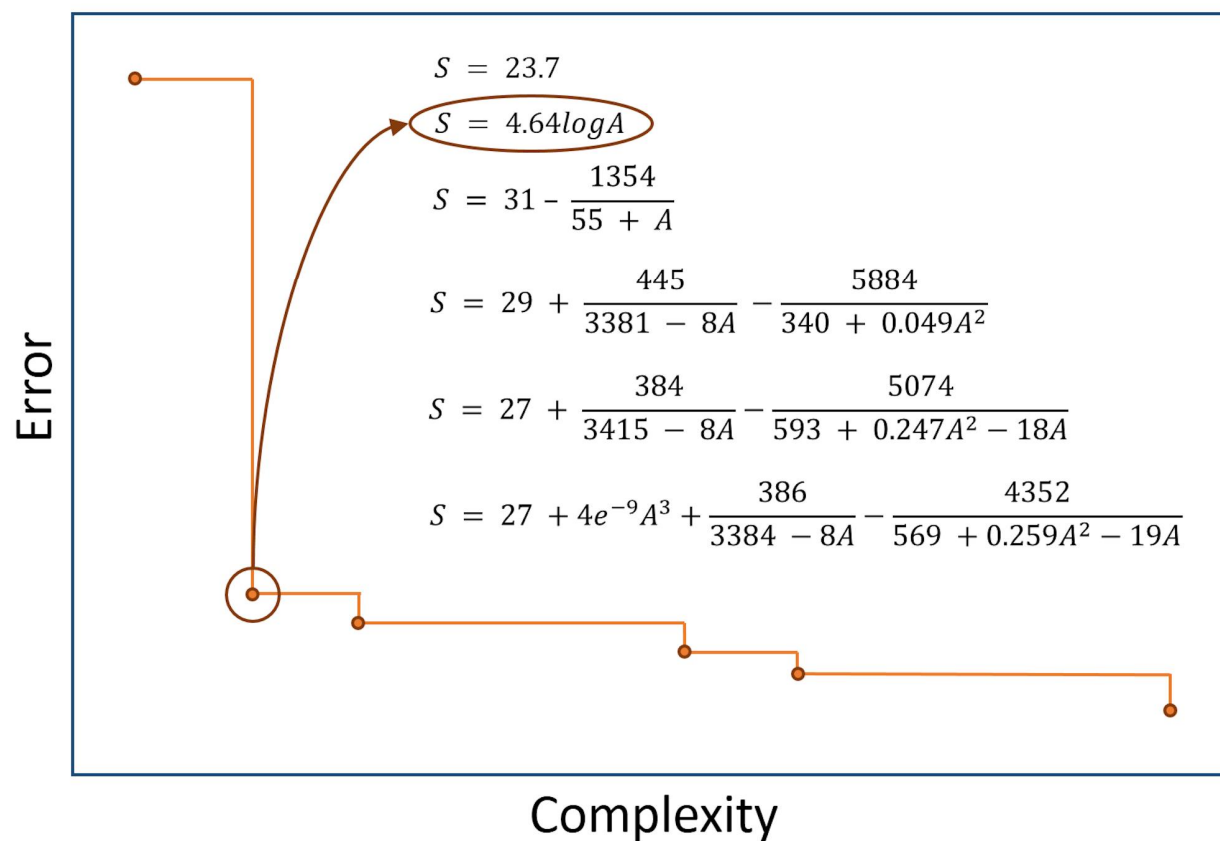
629

32

630 **Fig. 4. Accumulation curve for spider sampling in Gerês (Portugal).**

631 The result of searching for the best fitting asymptotic formula using symbolic regression is also

632 shown.

633

634

635

636    **Figure 1**

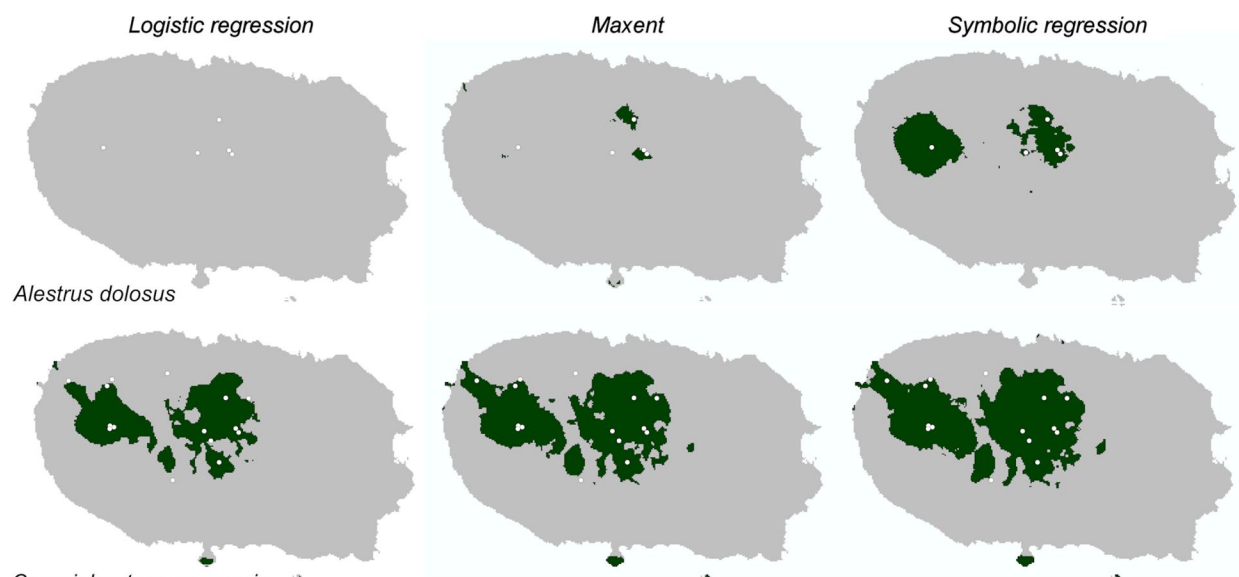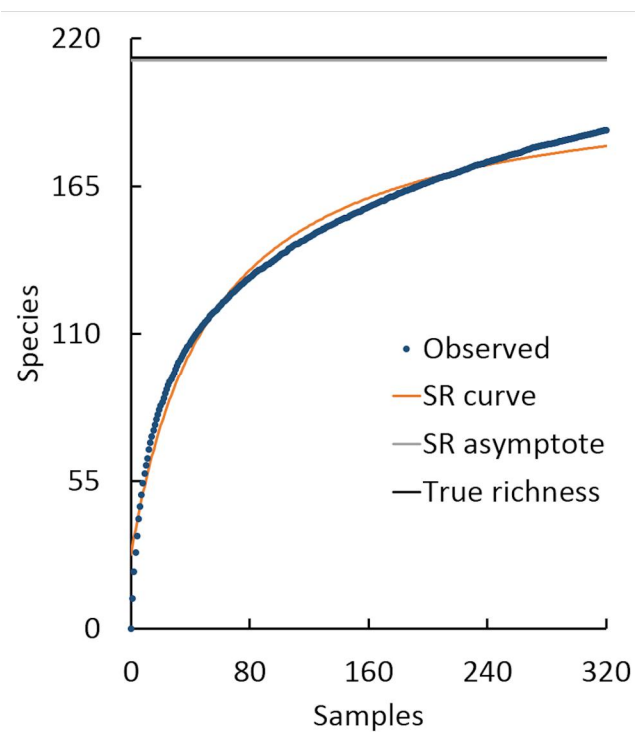$$S = 23.7$$

$$S = 4.64 \log A$$

$$S = 31 - \frac{1354}{55 + A}$$

$$S = 29 + \frac{445}{3381 - 8A} - \frac{5884}{340 + 0.049A^2}$$

$$S = 27 + \frac{384}{3415 - 8A} - \frac{5074}{593 + 0.247A^2 - 18A}$$

$$S = 27 + 4e^{-9}A^3 + \frac{386}{3384 - 8A} - \frac{4352}{569 + 0.259A^2 - 19A}$$

637

638    **Figure 2**

639

640 **Figure 3**

641

642 **Figure 4**

643 **Appendix 1. Examples of general principles in ecology and of some of the respective**

644     **statistical models.**

645

646 **Appendix 2. Data and settings used for all analyses in the paper (Eureqa file:**

647     **http://www.nutonian.com/products/eureqa/).**

648

649 **Appendix 3. Species, area and age for each Azorean or Canarian Island.**