# Species Profiles Support Recommendations for Quality Filtering of Opportunistic Citizen Science Data

Camille Van Eupen*,[a,b], Dirk Maes[c,d], Marc Herremans[e], Kristijn R.R. Swinnen[e], Ben Somers**,[b] and Stijn Luca**,[a]

* corresponding author

** joint last authors

[a] Ghent University, Department of Data Analysis and Mathematical Modelling, Coupure Links 653, B-9000 Ghent, Belgium; camille.vaneupen@kuleuven.be (ORCID 0000-0002-0924-8892); stijn.luca@ugent.be (ORCID 0000-0002-6781-7870)

[b] KU Leuven, Department of Earth and Environmental Sciences, Division Forest Nature and Landscape, Celestijnenlaan 200E, B-3001 Heverlee, Belgium; ben.somers@kuleuven.be (ORCID 0000-0002-7875-107X)

[c] Research Institute for Nature and Forest (INBO), Herman Teirlinckgebouw, Havenlaan 88 box 73, B-1000 Brussels, Belgium; dirk.maes@inbo.be (ORCID 0000-0002-7947-3788)

[d] Radboud Institute for Biological and Environmental Sciences (RIBES), Radboud University, PO Box 9010, NL-6500 GL Nijmegen, The Netherlands; dirk.maes@ru.nl

[e] Natuurpunt Studie, Coxiestraat 11, 2800 Mechelen, Belgium; ; marc.herremans@natuurpunt.be (ORCID: 0000-0002-9719-8732); kristijn.swinnen@natuurpunt.be (ORICD 0000-0002-1910-9247)

1    **ABSTRACT**

2    Opportunistic citizen science data are commonly filtered in an attempt to improve their applicability for

3    relating species occurrences with environmental variables. Recommendations on when and how to filter,

4    however, have remained relatively general and associations between species traits and filtering

5    recommendations are sparse. We collected six traits (body size, detectability, classification error rate,

6    familiarity, reporting probability and range size) of 52 birds, 25 butterflies and 14 dragonflies. Both

7    absolute (values not rescaled) and relative traits (values rescaled per taxonomic group) were linked to

8    filter effects, i.e. the impact on three different measures of species distribution model performance

9    caused by applying three different quality filters, for different degrees of sample size reduction. First, we

10   applied multiple regressions that predicted the filter effects by either absolute (including taxonomic

11   group) or relative traits. Second, a principal component and clustering analysis were performed to define

12   five species profiles based on species traits that were retained after a multiple regression model selection.

13   The analysis of the profiles indicated the relative importance of species traits and revealed new insights

14   into the association of species traits with changes in model performance after data quality filtering. Both

15   taxonomic group (more than absolute traits) and relative species traits (mainly classification error rate,

16   range size and familiarity) defined the impact of data quality filtering on model performance and we

17   discourage the selection of a quality filtering strategy based on one single species trait. Results further

18   confirmed the importance of considering the goal of the study (i.e. increasing model discrimination

19   capacity, sensitivity or specificity) as well as the change in sample size caused by stringent filtering. The

20   general species knowledge among citizen scientists (importance of observer experience), together with

21   the mechanism of record verification in an opportunistic data platform (importance of verifiable

22   metadata) have the largest potential for enhancing the quality of opportunistic records.

23   **KEYWORDS**

24    data quality filtering, filtering recommendations, opportunistic data, presence-only data, species

25    distribution models, species traits

26 **1. Introduction**

27 Biodiversity conservation needs adequate monitoring of species (Lindenmayer et al., 2020), especially in

28 times of rapid changes to the environment that threaten species and reduce abundances at alarming rates

29 (Newbold et al., 2015; Urban et al., 2016). Structured surveys, where species are recorded in a

30 standardized manner, are commonly put forward as the most desirable strategy for biodiversity

31 monitoring because of their high information content (Dobson et al., 2020). It is challenging, however, to

32 organize structured surveys for a large variety of species over broad spatial and temporal scales, leading

33 to spatial and temporal data gaps (Urban et al., 2016). When information on the potential distributions of

34 species is needed, e.g. for assigning protected areas (Thomaes et al., 2008) or areas of high potential

35 nature value (Maes et al., 2005), spatial gaps can be filled by using species distribution models (SDMs).

36 These models link environmental parameters, such as landscape and climate variables, to species

37 occurrence records (Guisan and Zimmermann, 2000).

38 Species occurrence records with high information content (e.g. collected in structured surveys) are

39 preferably used as input for SDMs, but as aforementioned, such data are sparse. Therefore, SDMs are

40 increasingly built with data collected by citizen scientists in a semi-structured manner (termed semi-

41 structured data) or in an opportunistic unstructured manner (termed opportunistic data). Large citizen

42 science initiatives that have online data platforms either focus on one data type (e.g. eBird contains only

43 semi-structured data (Sullivan et al., 2009), iNaturalist contains only opportunistic data

44 (https://www.inaturalist.org/)) or they combine both data types (e.g. iRecord

45 (https://www.brc.ac.uk/irecord/), waarnemingen.be (https://www.waarnemingen.be) and

46 Observation.org (https://observation.org)). The majority of opportunistic data consist of species presence

47 records with some basic information such as the date and geographical precision of the observation

48 (termed opportunistic presence-only data). The main advantage of opportunistic data is the availability in

49 large amounts, over large geographical areas and potentially long periods (Kosmala et al., 2016). A major

50    disadvantage, however, is the prevalence of different types of bias and error (Bird et al., 2014; Isaac and

51    Pocock, 2015) caused by a lack of standardised design and accompanying metadata. Opportunistic citizen

52    science data is therefore associated with uncertainty and scepticism towards its use for SDMs (Burgess et

53    al., 2017), and has led to many publications on how and when to use this type of data for biodiversity

54    research (e.g. Henckel et al., 2020; Isaac & Pocock, 2015; Maes et al., 2015; Van Strien et al., 2013). Ideally,

55    the benefits of both unstructured opportunistic data (quantity) and (semi-)structured survey data

56    (information content) are exploited simultaneously. This can be in the form of model-based data

57    integration (for a recent review see Isaac et al., 2020) or when structured data is used as external

58    validation data for SDMs (e.g. Matutini et al., 2021; Van Eupen et al., 2021).

59    In preparation for an SDM study, data are expected to be cleansed (Zurell et al., 2020). At a minimum, this

60    includes the removal of spatial and temporal outliers, duplicates and records with low precision (Serra-

61    Diaz et al., 2017). When dealing with opportunistic data, cleansing usually also implies stringent filtering,

62    where data are filtered based on record attributes that hold information on the observation process or

63    post-entry data validation (Steen et al., 2019). Even though drawing direct ecological inferences from

64    opportunistic observations is not recommended (Dobson et al., 2020), clear recommendations on when

65    and how to filter opportunistic data remain sparse (but see e.g. Kamp et al., 2016; Steen et al., 2019; Van

66    Eupen et al., 2021; Vantieghem et al., 2017). The study of Van Eupen et al. (2021) highlighted the

67    importance of considering both the type of filter and the resulting change in sample size, yet variation

68    among species in their response to data quality filtering remained large. Grouping species according to

69    their taxonomy revealed that filtering benefitted some groups (i.e. plants and dragonflies) more than

70    others (i.e. butterflies and birds). In this paper, we aim to verify whether grouping species according to a

71    priori selected life-history and/or ecological traits could better substantiate recommendations for data

72    quality filtering.

73 Species traits have been linked extensively to SDM performance and those that cause most variation can

74 usually (but not exhaustively) be compiled to the following three: (1) traits that define the species-

75 environment relationship (e.g. range size, niche breadth (Brotons et al., 2007; Stockwell and Peterson,

76 2002) and habitat association (Chefaoui et al., 2011)), (2) traits that impact the detectability of the species

77 in space and time (e.g. conspicuousness (Seoane et al., 2005), migratory behaviour (Carrascal et al., 2006)

78 and lifespan (Hanspach et al., 2010)), and (3) traits that influence the proneness to misidentification (e.g.

79 phylogenetic relatedness (Vantieghem et al., 2017)).

80 Notwithstanding the vast amount of proof on the link between species traits and absolute SDM

81 performance, few studies have successfully linked species traits to the change in SDM performance caused

82 by stringent filtering of species occurrence records (but see e.g. Steen et al., 2019, where models of more

83 restricted species performed better when using data collected with lower effort). This could be due to the

84 higher quality of the unfiltered data in most of these studies (e.g. semi-structured data in Steen et al.

85 (2019)) or due to the conflicting character of the simultaneous impact of data quality filtering, i.e. an

86 increase in data quality and a decrease in sample size (Van Eupen et al., 2021). By assessing this twofold

87 effect on an extensive dataset of opportunistic records, *waarnemingen.be*, we will aid the optimisation of

88 the data cleansing process that is essential for high-quality SDMs (Zurell et al., 2020).

89    **2.    Material and methods**

90       **2.1. Species data and impact of quality filtering**

91    We used a dataset from a previous study on data quality filtering (Van Eupen et al., 2021, Van Eupen et

92    al., 2021b), where three dichotomous filters were applied to opportunistic species observations belonging

93    to four well-studied taxonomic groups in Flanders, i.e. birds, butterflies, dragonflies and plants. Plant

94    observations were not used in the present analysis because their traits are not directly comparable to

95    animal species traits. Data were collected from the '*waarnemingen.be*' database, an online citizen science

96    platform that contains both semi-structured and opportunistic records (Swinnen et al., n.d.). The selected

97    dataset included year-round records from January 2014 to September 2019 with a geographical precision

98    of 500 metres or smaller (records can be submitted to the platform as point locations with specified

99    precision or as observations made within a larger area). Potential seasonal changes in habitat occupancy

100   or range size were not considered. The selection excluded records verified as incorrect by species experts

101   or the data platform's auto-validation system (Swinnen et al., 2018) and records from non-native or non-

102   breeding birds in Flanders (Vermeersch et al., 2020). Opportunistic presence-only data were used for

103   model training, also excluding absences (zero-counts), and semi-structured data for model testing (Van

104   Eupen et al., 2021). The three filters were: '**ACTIVITY**', based on an observer's average annual activity rate,

105   where the filter consists in removing records from less active observers; '**DETAIL**', based on the presence

106   of metadata beyond default requirements (i.e. species name, location, date and observer id), where the

107   filter consists in removing records that were submitted without any additional information (e.g. sex,

108   count, behaviour); and '**VALSTAT**', based on the validation status of a record in the data platform, where

109   the filter consists in removing doubtful and unevaluated records (Table 1). These are all records that could

110   not be verified by species experts because key information was missing or because the record was not

111   assessed yet by an expert at the moment the dataset was extracted.

112 *Table 1: Overview and definitions of the used variables in this study.*

| Data quality filters | description: | based on: |
|---|---|---|
| ACTIVITY | removes records from less active observers | an observer's average annual activity rate |
| DETAIL | removes records that were submitted without any additional information | the presence of metadata beyond default requirements |
| VALSTAT | removes doubtful and unevaluated records | the validation status of a record in the data platform |
| **Species traits** | *description:* | *source:* |
| Body size | wing length (birds and butterflies) or head-to-tail length (dragonflies) | Bink (1992); Storchová and Hořák (2018), https://www.vlinderstichting.nl/libellen/ |
| Classification error rate | the number of erroneous photo records (i.e. observations accompanied by a photograph) relative to the total number of photo records. | the *waarnemingen.be* data portal during the study period |
| Detectability | the probability of detecting a species on the condition that it is present | quantified by applying site occupancy models to complete checklist data, retrieved from the *waarnemingen.be* data portal |
| Familiarity | reflects how well-known a species is by the average observer | number of search results retrieved from the Google search engine |
| Reporting probability | the likelihood that a species is reported by an average observer, on the condition that it is present and that the taxonomic group it belongs to is surveyed | a species' relative (per taxonomic group) average reporting rate divided by its detectability, retrieved from the *waarnemingen.be* data portal |
| Range size | the distribution range size | the total number of grid cells (km²) in which a species has been recorded during the study period, retrieved from the *waarnemingen.be* data portal |
| Absolute traits | unscaled trait values as retrieved by the different methods described | |
| Relative traits | scaled trait values; using the following transformation per taxonomic group: $$y = \frac{x - \min(x)}{\max(x) - \min(x)}$$ | |
| **Impact on model performance** | | |
| Δ AUC | change in the area under the receiver operating characteristic | |
| Δ sensitivity | change in the true positive rate (TPR) after data quality filtering $$TPR = \frac{true\ positives}{true\ positives + false\ negatives}$$ | |
| Δ specificity | change in the true negative rate (TNR) after data quality filtering $$TNR = \frac{true\ negatives}{true\ negatives + false\ positives}$$ | |
| **Filter effects** | | |
| All combinations of data quality filters and impact on model performance | the impact of data quality filtering (Δ AUC, Δ sensitivity and Δ specificity) by the three filters ACTIVITY, DETAIL and VALSTAT | |
| **Sample size situations** | | |
| real | the actual reduction in the number of presences after data quality filtering | |
| r50 | a relative reduction in the number of presences after data quality filtering of more than 50% | |
| ss100 | a reduction to 100 presences after data quality filtering | |

113    After filtering, records were aggregated to a 1x1 km resolution to reduce spatial bias (Kramer-Schadt et

114    al., 2013). Species were selected based on opportunistic data availability for model training as well as

115    semi-structured data availability for model validation (Van Eupen et al., 2021). The impact of filtering on

116    the performance of Maxent (Phillips et al., 2006) was assessed by evaluating the difference in three

117    commonly used evaluation metrics of model discrimination before and after filtering: the area under the

118    receiver operating characteristic (AUC), sensitivity and specificity (Fielding and Bell, 1997). Their use was

119    justified because models were run for the same geographical extent and model predictions were

120    evaluated on a testing set that contained the same presences and absences per species (Jiménez-

121    Valverde, 2012; Lobo et al., 2008). Model performance could therefore be interpreted in a relative

122    manner, where an increase in AUC after filtering implied that the filtered data produced models that could

123    better distinguish between testing presences and absences, and increases in sensitivity and specificity

124    implied a higher predicted positive and negative fraction respectively. For the analysis in this study, we

125    extracted the change in model performance (i.e. Δ AUC, Δ sensitivity and Δ specificity) (Table 1), after

126    using the three single filters (ACTIVITY, DETAIL and VALSTAT) for 52 birds, 25 butterflies and 14

127    dragonflies. For a summary per species of the data used for model testing and model training (unfiltered

128    and filtered data) and of the impact on model performance, we refer to the supplementary information

129    1 (Table C.1) and supplementary information 2 respectively in the study of Van Eupen et al. (2021).

### 2.2.  Species traits

131    We used six species traits that can be related to data quality in opportunistic citizen science data based

132    on literature review and expert opinion: body size, detectability, classification error rate, familiarity,

133    reporting probability and range size (Table 1). Abundance was not considered because the largely

134    unstructured *waarnemingen.be* database contains unreliable count data that are mostly without a clear

135    reference of time and space. All trait values can be found in the supplementary information (Table S1).

136 **Body size** equals the wing length for birds (Storchová and Hořák, 2018) and butterflies (Bink, 1992) and

137 head-to-tail length for dragonflies (https://www.vlinderstichting.nl/libellen/).

138 The **classification error rate** reflects how likely it is for an average observer to wrongly identify a species.

139 This was quantified by the number of erroneous photo records (i.e. observations accompanied by a

140 photograph) of a species in the *waarnemingen.be* data portal, relative to its total number of photo

141 records. The portal keeps track of changes in the identification of a species, and we considered only the

142 changes at the species-level as erroneous (and for example not the changes from family or genus to

143 species-level). Auto-corrections made by the observer were excluded.

144 **Detectability** is the probability of detecting a species on the condition that it is present (MacKenzie et al.,

145 2017). Species detectability was retrieved from applying site occupancy models to complete semi-

146 structured checklist data extracted from *waarnemingen.be*, following Johnston et al. (2021). Detection

147 histories consisted of five to ten repeated visits to a specific site (a 1 km grid cell) by the same observer in

148 a period of closure (i.e. a period with no supposed changes in occupancy). A period of closure was defined

149 as 20 consecutive days in the peak active season of a species. The peak active season was defined as every

150 10 days with an observation count above the average count of all observations in a year, excluding egg,

151 larva, pupa and caterpillar observations. Covariates used to describe the detection process were: checklist

152 duration (in minutes), starting time of the checklist, search effort (i.e. the number of species recorded at

153 a specific location, supporting on the principle of species accumulation curves (Colwell et al., 2004)), and

154 open habitat (grasslands, wetland, marshes and water) versus closed habitat (forest and woodland),

155 because of an increased detectability (visually and, for birds, also auditory) in open habitat types (Johnston

156 et al., 2014; Morton, 1975). Detection probabilities were predicted for all grids with covariate values and

157 averaged to attain one value per species.

158    **Familiarity** refers to how well-known a species is by the average observer and was quantified by the

159    number of Belgian websites with the Dutch name of the species in the title, retrieved from the Google

160    search engine (Żmihorski et al., 2013). We added two extra search terms that specified the taxonomic

161    group (in Dutch) and excluded the *waarnemingen.be* website to avoid counting individual observations

162    on the used data platform, e.g. *"Bruinrode Heidelibel" site:.be libel -waarnemingen.be.* An Incognito

163    window was used to unlink search results from the used google account.

164    **Range size** is the distribution range size of the species during the entire study period 2014-2019 in the

165    study area and was quantified as the total number of grid cells (km²) in which a species has been recorded
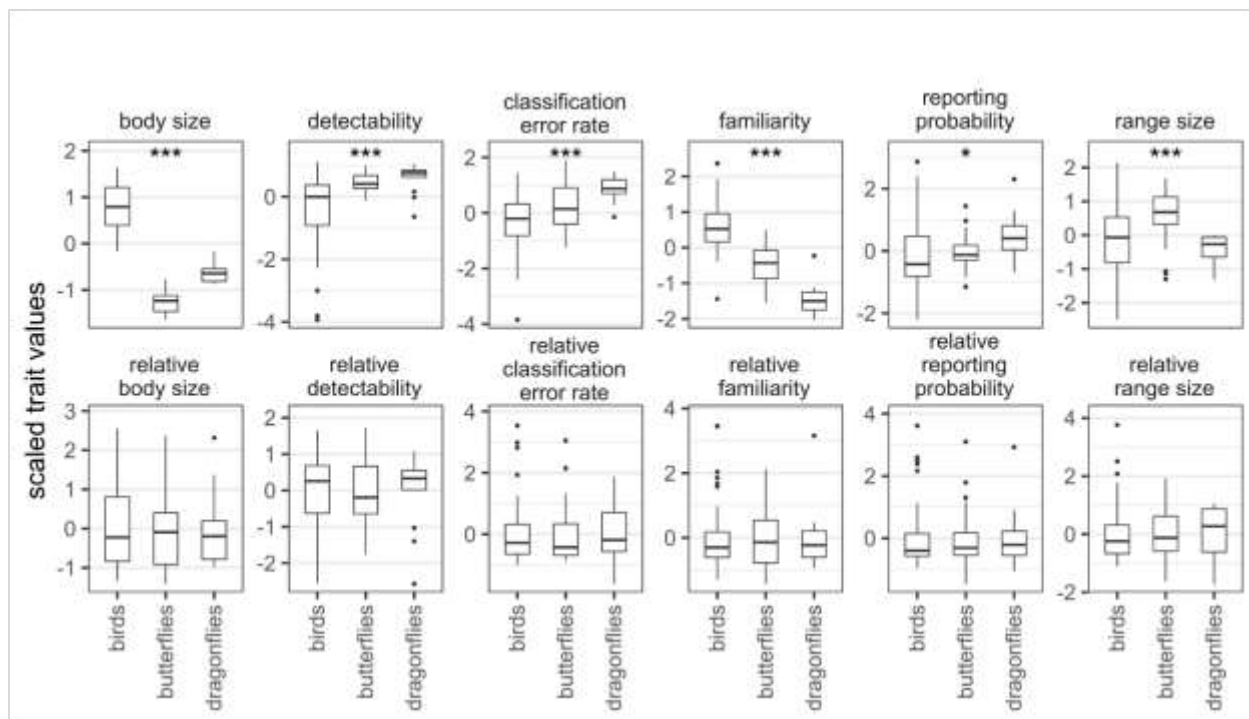
166    (McPherson et al., 2004).

167    **Reporting probability** is the likelihood that a species is reported by an average observer, on the condition

168    that it is present and that the taxonomic group it belongs to is surveyed. To meet these requirements, we

169    looked at the peak of the active season and calculated the relative number of species observations to the

170    number of observations for a taxonomic group. This was averaged across locations and observers. We

171    subsequently divided this number by the average detectability across locations where the species was

172    present to correct for the impact of detectability on reporting rate.

173    **2.3.  The impact of data quality filtering**

174    To build recommendations for data quality filtering based on species traits, we first analysed the

175    multivariate relationship between species traits and the filter effects. Consequently, species were

176    grouped in species profiles characterised by the most highly associated traits to assess if such groups

177    presented a similar response to data quality filtering. By filter effect, we mean the impact of data quality

178    filtering by the three filters ACTIVITY (only observations from active observers), DETAIL (only detailed

179    observations) and VALSTAT (only approved observations) on three evaluation metrics: AUC, sensitivity

180    and specificity. All analyses were conducted in R (R Core Team, 2021).

181 ### 2.3.1. Multi-trait analysis

182 Relationships between species traits and filter effects were examined using multiple (multi-trait)

183 regressions. The data were modelled in beta-regressions (*betareg* package v3.1-4, Cribari-Neto & Zeileis,

184 2010), because of the bound character of the response variable (Δ AUC, Δ sensitivity and Δ specificity

185 theoretically range from -1 to 1). Filter effect values were rescaled to fall between 0 and 1 with the

186 following transformation: $y = \frac{x - \min(x)}{\max(x) - \min(x)}$. To reduce the impact of outliers, data points with a cook's

187 distance of more than four times the mean cook's distance of all data points were removed (Ferrari et al.,

188 2004). As trait values showed taxonomic differences (Figure 1), continuous values (absolute traits) were

189 rescaled per taxonomic group (relative traits), using the aforementioned transformation (Table 1). The

190 relative values can be informative for patterns across taxonomic groups that would go unnoticed

191 otherwise (e.g. birds are always larger than butterflies, but similar impacts from filtering might be

192 observed for large birds as well as large butterflies).

193

*Figure 1: Summary of the species traits per taxonomic group after value transformation and standardisation. Absolute traits (top*

*row) were rescaled to relative traits (bottom row) per taxonomic group to assess patterns across taxonomic groups. Stars indicate*

*differences in the medians of the trait values between taxonomic groups (\*\*\* = p < 0.001, \* = p < 0.05).*

197 First, multi-trait regressions were performed using the log-transformed **absolute** trait values as

198 continuous variables and the taxonomic group as a factor variable. Second, **relative** traits were regressed

199 against the filter effects. Trait values were standardised and multicollinearity was reduced by retaining

200 only those variables with a Variance Inflation Factor (VIF) below 5 (Menard, 2001). We modelled the

201 absolute and relative traits separately because of high pairwise correlations among most of these

202 variables (Figure S1). We also quantified variable importance by leaving out each trait one by one and

203 calculating the decrease in pseudo-$R^2$ compared to the full model. Finally, we performed a model selection

204 based on three conditions to obtain parsimonious models for each filter effect: (1) the increase in the

205 Akaike's Information Criterion (AIC) had to be smaller than (a conservative) five (Burnham et al., 2011),

206 (2) the model should at least contain the most important variable and (3) the simplest model was selected

207 (i.e. the model with the least parameters).

### 208    2.3.2.    Species profiles

209    To test whether groups of species with similar traits can improve recommendations for data quality

210    filtering, we delineated species profiles. Species were clustered into groups with similar traits using the

211    FactoMiner package v1.34 (Le et al., 2008). This package performs a principal component analysis (PCA)

212    on a set of variables (i.e. species traits) followed by an agglomerative hierarchical clustering of individuals

213    (i.e. species) each described by principal components of those variables (Husson et al., 2010). The active

214    variables, included in the PCA and clustering, were those variables that contributed most to the change in

215    model performance across filters, resulting from the multi-trait regression model selection.

216    Supplementary variables were added to characterise the clusters further, without impacting the clustering

217    itself, and comprised: the remaining traits and the impact of filtering on model performance per filter

218    (quantitative), and the taxonomic group (qualitative).

219    We delineated the profiles based on four conditions. First, traits used were those remaining after model

220    selection (see section 2.3.1) in at least one multi-trait regression. Second, the number of clusters was

221    chosen based on the increase of inertia between two consecutive aggregation steps in the hierarchical

222    tree (Husson et al., 2010). Third, profiles were ideally associated with one or more distinctive filter effects:

223    (1) an increase in AUC, (2) a decrease in AUC, (3) an increase in sensitivity and/or decrease in specificity

224    or (4) an increase in specificity and/or decrease in sensitivity. Fourth, profiles should be ecologically

225    meaningful, where we relied on species experts to evaluate the profiles' species composition. To this end,

226    we experimentally changed the number of clusters before selecting the final profiles, by choosing

227    different heights in the hierarchical tree.

### 228    2.3.3.    Impact of sample size

229    The role of sample size in the relationship between species traits and filter effects was assessed by adding

230    two sample size situations based on previous recommendations (Van Eupen et al., 2021), where filtering

231    was not advised when sample size was reduced by more than 50% or when the resulting sample size was

232    100 presences. We looked at (1) the **real** situation, i.e. the filter effects when sample size after filtering

233    was not altered (sample size was reduced by an amount that depended on the applied filter), (2) the **r50**

234    situation, i.e. the filter effects when sample size was reduced by 50% or more and (3) the **ss100** situation,

235    i.e. the filter effects when sample size was reduced to 100 presences. Adding these two situations could

236    aid interpretation and simulate situations occurring in datasets of lower quality (i.e. where fewer

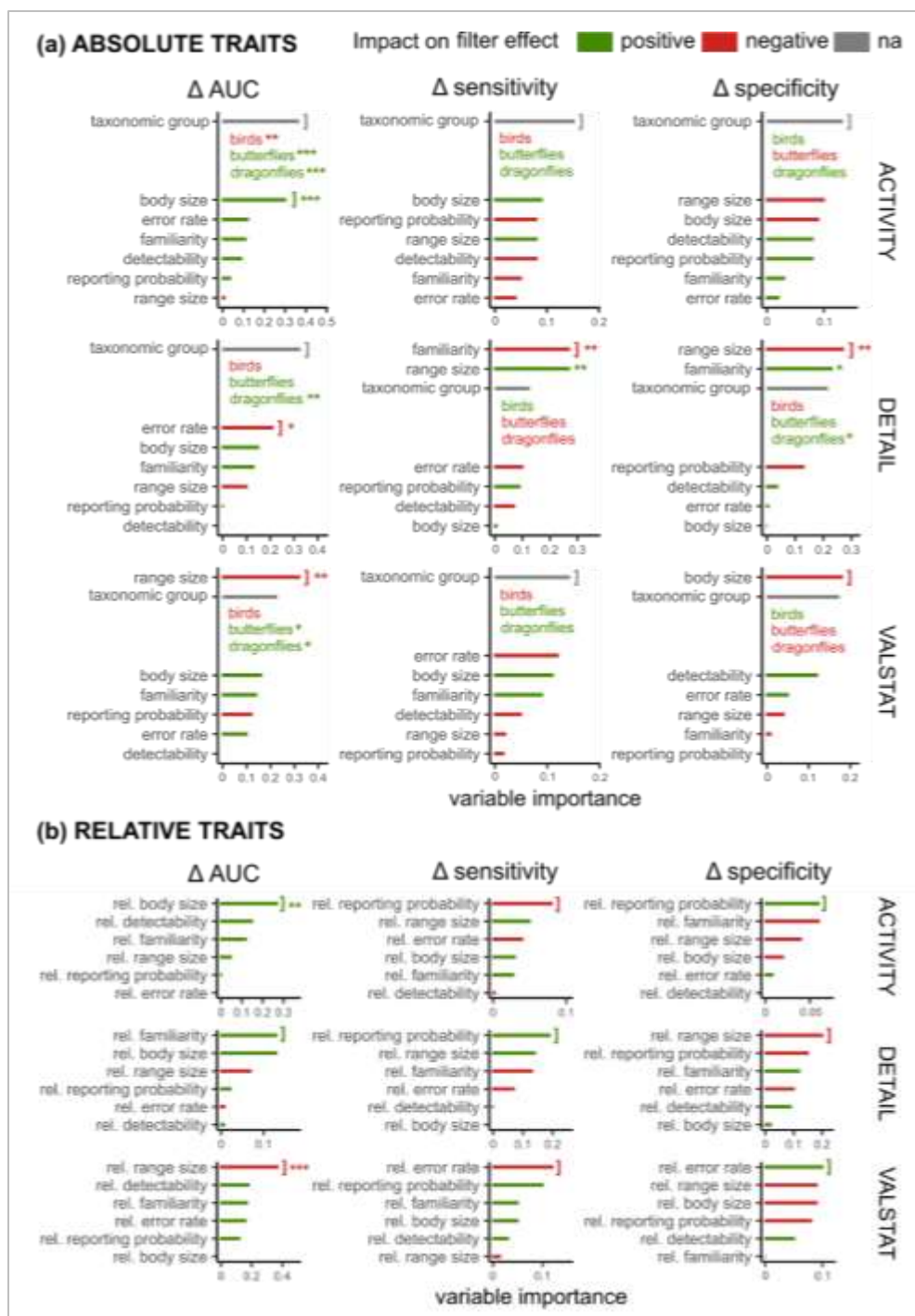237    presences are kept after stringent filtering).

238 **3. Results**

239 **3.1. Multi-trait analysis**

240 Figure 2 shows that the relative importance of the different traits in their association with the filter effects

241 varies among filters and model evaluation metrics. Considering the absolute traits (Figure 2a), the

242 taxonomic group was the most important variable in five out of nine cases. When the goal was to increase

243 AUC, it was best to use data from active observers or approved observations for butterflies and

244 dragonflies, or detailed observations for dragonflies. Specificity could be increased for dragonflies by using

245 detailed observations. No other significant differences ($p < 0.05$) between taxonomic groups were

246 detected in the multiple regressions. AUC of models from large-bodied species could best be increased by

247 using data from active observers, AUC of models from species with a low error rate by using detailed

248 observations and AUC of models from species with a restricted range size by using approved observations.

249 Sensitivity of models from unfamiliar species benefitted from using detailed observations, as did

250 specificity of models from species with restricted range sizes. Specificity of models from small-bodied

251 species benefitted from using approved observations.

252 Considering the relative traits (Figure 2b), using observations from active observers worked best for large-

253 bodied species (to increase AUC), for species with low reporting probability (to increase sensitivity) or for

254 species with high reporting probability (to increase specificity). Using detailed observations most

255 benefitted familiar species (to increase AUC), species with high reporting probability (to increase

256 sensitivity) or species with a restricted range size (to increase specificity). Using approved observations

257 was most valuable for species with a restricted range size (to increase AUC), for species with a low

258 classification error rate (to increase sensitivity) or for species with a high classification error rate (to

259 increase specificity).

260    Multicollinearity among absolute and relative traits was negligible (VIF < 5), so all traits were included in

261    the multi-trait regressions. Neither absolute nor relative detectability was retained after model selection

262    as these traits explained less variation compared to others. We did observe that detectability was

263    negatively correlated with familiarity (r = -0.38, p < 0.001), which can be explained by the taxonomic

264    differences found in both traits (Figure 1 and Figure S1). Detectability was also negatively correlated with

265    reporting probability (r = -0.62 and r = -0.65 for absolute and relative traits respectively, p < 0.001), which

266    can be explained by their inverse dependence (Figure S1 and section 2.2).

Figure 2: Variable importance in the multi-trait regressions for absolute (a) and relative (b) species traits per filter (ACTIVITY,

DETAIL and VALSTAT) and change in model evaluation metric (Δ AUC, Δ sensitivity, Δ specificity). Variable importance is expressed

as the square root of the change in pseudo-$R^2$ when leaving out one variable at a time from the full model. Colours indicate a

positive (green) or negative (red) impact of the trait on the filter effect, factor variables have grey (n/a) colour. Square brackets

indicate the variables kept after model selection (i.e. the simplest model with an increase in the Akaike's Information Criterion

273    *(AIC) of less than 5 compared to the best model where at least the most important variable was included). Asterisks indicate*

274    *significant model coefficients (\*\*\* = p < 0.001, \*\* = p < 0.01, \* = p < 0.05).*

275    **3.2.  Species profiles**

276    The active variables that we used in the PCA and clustering analysis were the continuous variables kept

277    after model selection in the multi-trait analysis: i.e. body size, relative body size, classification error rate,

278    relative classification error rate, familiarity, relative familiarity, relative reporting probability, range size

279    and relative range size. In the experimental phase, three clusters best captured the variation in species

280    traits, but the impact on model performance still showed large variation within profiles. Ecologically, these

281    profiles also separated species into quite general groups and we cross-checked the clustering of the

282    species into four or more different profiles with species experts. Five clusters appeared the best outcome

283    while keeping cluster size at a reasonable level (minimum cluster size equalled 7 species) (Table 2). The

284    full results of the PCA and clustering analysis are presented in the supplementary information (Table S2

285    and Figure S2).

286    Positive or negative recommendations were only noted when the goal was to increase AUC. When the

287    goal was to increase sensitivity or specificity, recommendations were either cautious or alarming and in

288    most cases, filter recommendations for increasing sensitivity and specificity were opposite to each other.

289    Similar impacts between profiles on one evaluation metric might have a different impact on other metrics

290    (e.g. similar impact on AUC but a different impact on sensitivity and specificity in profiles 1 and 4).

291    *Table 2: Recommendations for data quality filtering for the five species profiles, described by five relative traits (body size,*

292    *classification error rate, familiarity, reporting probability and range size) and four absolute traits (body size, classification error*

293    *rate, familiarity and range size). Recommendations are positive (green – all values in the 90% confidence interval are positive),*

294    *cautious (blue - the average filter effect is positive but the 90% confidence interval also includes negative values), alarming (orange*

295    *- the average filter effect is negative but the 90% confidence interval also includes positive values) or negative (red - all values in*

296    *the 90% confidence interval are negative). The taxonomic distribution of the species is given, as well as the most characterising*

297    *species per profile (in bold are the species closest to the cluster centre and in italic are the species furthest away from the other*

298    *cluster centres). The asterisks indicate the significance level at which traits, filter effects or taxonomic groups are associated with*

299    *a profile (\*\*\* = 0.001, \*\* = 0.01, \* = 0.05). For taxonomic groups, (+) and (-) indicate whether the group is significantly more or*
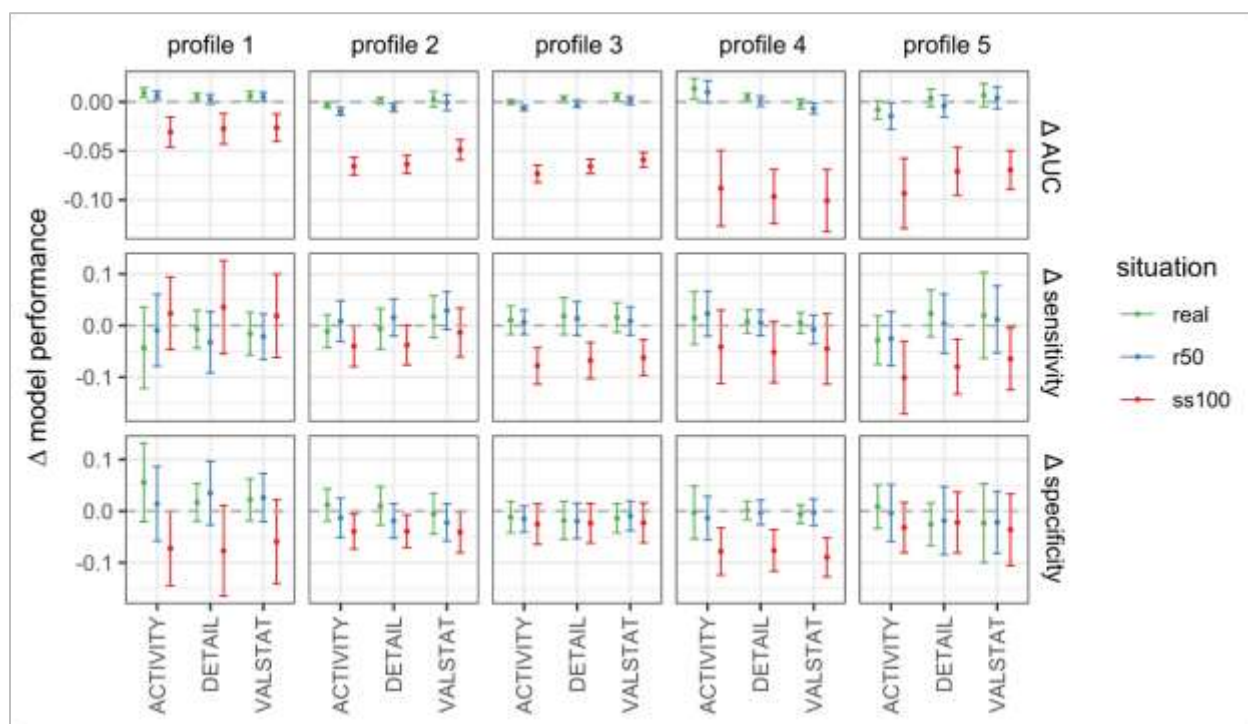
300    *less represented in a profile.*

| | PROFILE 1 | PROFILE 2 | PROFILE 3 | PROFILE 4 | PROFILE 5 |
|---|---|---|---|---|---|
| *Relative Traits* | High error rate *** <br><br> Widespread *** | Small body size *** <br><br> Restricted range size ** <br><br> Low error rate * <br><br> Unfamiliar * | Large body size *** <br><br> Restricted range size ** <br><br> Unfamiliar * | Familiar *** <br><br> Widespread *** <br><br> Large body size ** <br><br> Low error rate * | Familiar *** <br><br> High reporting probability *** <br><br> Low error rate * |
| *AUC* <br><br> *recommendations* | ACTIVITY ** > VALSTAT > DETAIL | VALSTAT > DETAIL <br><br> ACTIVITY ** | VALSTAT <br><br> DETAIL <br><br> ACTIVITY | ACTIVITY ** > DETAIL <br><br> VALSTAT | VALSTAT > DETAIL <br><br> ACTIVITY * |
| *sensitivity* <br><br> *recommendations* | ACTIVITY < VALSTAT < DETAIL | VALSTAT <br><br> ACTIVITY < DETAIL | DETAIL > VALSTAT > ACTIVITY | ACTIVITY > DETAIL > VALSTAT | DETAIL > VALSTAT <br><br> ACTIVITY |
| *specificity* <br><br> *recommendations* | ACTIVITY > VALSTAT > DETAIL | ACTIVITY > DETAIL <br><br> VALSTAT | DETAIL < VALSTAT < ACTIVITY | DETAIL <br><br> VALSTAT < ACTIVITY | ACTIVITY <br><br> DETAIL < VALSTAT |
| *Taxonomic group* | 20 species <br><br> 4 birds *** (-) <br><br> 6 butterflies <br><br> 10 dragonflies *** (+) | 35 species <br><br> 20 birds <br><br> 11 butterfly <br><br> 4 dragonflies | 17 species <br><br> 17 birds *** (+) | 12 species <br><br> 4 birds <br><br> 8 butterflies ** (+) | 7 species <br><br> 7 birds * (+) |
| *Absolute traits* | High error rate *** <br><br> Unfamiliar ** | Restricted range size *** <br><br> Small body size *** <br><br> Low error rate * | Large body size *** <br><br> Low error rate * | Widespread *** <br><br> Low error rate * | Familiar *** <br><br> Low error rate * |
| *Characterising species* | **Pieris napi** <br> ***Sympetrum striolatum*** <br> **Sympetrum sanguineum** <br> **Aeshna cyanea** <br> **Maniola jurtina** <br> *Pieris brassicae* <br> *Pieris rapae* <br> *Enallagma cyathigerum* <br> *Larus canus* | **Oenanthe oenanthe** <br> **Turdus pilaris** <br> **Tachybaptus ruficollis** <br> **Delichon urbicum** <br> **Rallus aquaticus** <br> *Platycnemis pennipes* <br> *Colias crocea* <br> *Calopteryx splendens* <br> *Pyrrhosoma nymphula* <br> *Motacilla alba* | **Tadorna tadorna** <br> **Circus aeruginosus** <br> **Numenius arquata** <br> **Egretta garzetta** <br> **Branta leucopsis** <br> *Cygnus olor* <br> *Branta canadensis* <br> *Anser anser* <br> *Ardea alba* <br> *Corvus frugilegus* | **Vanessa cardui** <br> **Polygonia c.album** <br> ***Gonepteryx rhamni*** <br> ***Vanessa atalanta*** <br> **Falco tinnunculus** <br> *Buteo buteo* <br> *Aglais io* <br> *Papilio machaon* | ***Cuculus canorus*** <br> ***Alcedo atthis*** <br> **Perdix perdix** <br> **Carduelis carduelis** <br> ***Athene noctua*** <br> *Ciconia ciconia* <br> *Alopochen aegyptiaca* |

301  Absolute traits appeared highly associated with the most represented taxonomic group for four out of

302  five profiles (Table 2 and Figure 1). Profile 1 contained more dragonflies, indeed species with a higher

303  classification error rate that are less familiar. Profile 4 contained more butterflies, species with a large

304  range size, yet not necessarily a lower error rate. Profile 3 contained birds only, which have larger body

305  sizes and lower error rates. There is a difference with profile 5 though, also containing only birds, where

306  higher familiarity and lower error rates are characterising traits. Profile 2 is not associated with one of the

307  three taxonomic groups, but species in this profile are mostly small, with a restricted range size and a

308  lower error rate, which are also relative traits that characterise this profile.

309  Recommendations based on relative traits were mostly similar to the results in the multi-trait analysis,

310  with a few exceptions (Table 2 and Figure 2). Model AUC for large species increased when using

311  observations from active observers, confirmed by negative and positive recommendations in profiles 2

312  and 4 respectively. In profile 3, however, body size seemed subordinate to the taxonomic group. Higher

313  reporting probability was associated with a higher Δ sensitivity when using detailed observations and with

314  a lower Δ sensitivity when using observations from active observers, confirmed in profile 5. Familiarity

315  had a positive impact on Δ AUC when using detailed observations (DETAIL), confirmed by positive and

316  cautious recommendations in profiles 4 and 5, yet only as the second-best option. Using DETAIL did not

317  necessarily worsen model AUC for unfamiliar species (profile 2 and 3), but not all species in these profiles

318  were unfamiliar (indicated by the weak significance level). For using only approved observations, range

319  size was a good indicator of a change in AUC (profiles 2, 3 and 4), except for the widespread species in

320  profile 1 where range size seemed to be subordinate to error rate. For using only detailed observations,

321  however, range size did not seem to drive filter recommendations when the goal was to increase model

322  specificity, except for the cautious recommendation for species with a restricted range size in profile 2.

323  Finally, the association between error rate and model sensitivity and specificity supported filter

324  recommendations when using only approved observations (profiles 1, 2, 4 and 5).

325    **3.3. Impact of sample size**

326    Reducing sample size further (beyond the already occurring decrease in sample size after data quality

327    filtering) impacted filter effects both positively and negatively (Figure 3). The impact on AUC was generally

328    negative, but the impact on sensitivity and specificity could in a few cases also be positive. It, therefore,

329    depends on the goal of the study whether reducing sample size further might have a desirable effect. Note

330    that there is usually a trade-off between sensitivity and specificity (when sensitivity increases, specificity

331    usually decreases and the other way around) (Jiménez-Valverde, 2012). Variability in the impact on model

332    performance also increased with decreasing sample size, except for species with a restricted range size

333    (profiles 2 and 3).

334


335    *Figure 3: Recommendations for data quality filtering for the five species profiles in the three sample size situations (real = actual*

336    *reduction in sample size when filtering, r50 = sample size reduced by at least 50%, ss100 = sample size reduced to 100 presences).*

337    *Dots and error bars are the means and 90% confidence intervals for the filter effects.*

338    Reducing sample size further mostly worsened model performance, especially when sample size was

339    reduced to 100 presences where recommendations became alarming or negative in most cases. In our

340    dataset, this meant that sample size was reduced by at least 77% (the lowest unfiltered sample size

341    equalled 432 presences). There were a few exceptions where reducing sample size further did have a

342    positive impact on model performance. For example, recommendations for increasing sensitivity could

343    change from alarming to cautious when reducing sample size over 50% for profile 2 (using observations

344    from active observers or detailed observations) and up till sample size reached 100 presences for profile

345    1 (all filters). Results also showed that model sensitivity was more (and specificity was less) impacted by

346    sample size reduction for profiles with birds only (profiles 3 and 5) compared to profiles with more

347    dragonflies and butterflies (profiles 1 and 4).

348    **3.4. Recommendations for data quality filtering**

349    Recommendations for data quality filtering were built on the various results presented in this article. In

350    general, users of opportunistic records should always pay attention when filtering reduces sample size by

351    more than half of its original size, leading to small sample sizes and we generally advise against filtering

352    when sample size is reduced by more than 75%. We further interpreted the filtering recommendations of

353    the PCA and clustering analysis (Table 2) together with the results of the multi-trait analysis (Figure 2). In

354    the following paragraphs, recommendations are formulated with the aim to increase AUC unless specified

355    otherwise.

356    Results showed that taxonomic group (more than absolute traits) and relative traits formed the best basis

357    for filtering recommendations and when we discuss traits in the following paragraphs, we mean the

358    relative values unless specified otherwise. We recommend using only data from active observers when

359    filtering opportunistic records of large or widespread butterfly and dragonfly species (profiles 1 and 4)

360    and approved observations when filtering bird records unless they are very familiar and widespread

361  (profile 4). In the cases where absolute traits were retained after model selection (Figure 2), it was the

362  relative rather than the absolute trait that was causing the filter effect. For example, dragonflies and

363  butterflies benefitted more from using observations from active observers (ACTIVITY) compared to birds,

364  yet a higher absolute body size also impacted this effect positively. This meant that dragonflies and

365  butterflies with a higher (relative) body size benefitted most from using the ACTIVITY filter. Keeping bird

366  observations from more active observers only was generally not recommended, except for widespread

367  species with a high classification error rate (profile 1).

368  Recommendations based on the taxonomic group seemed to overrule the impact of body size (profile 3)

369  and we advise against using body size as a motivator for filtering bird species data. Recommendations

370  based on the taxonomic group were also superior to the impact of familiarity (profiles 4 and 5), yet we

371  still recommend using more detailed observations (DETAIL) for familiar species, especially when they have

372  high reporting probability and an increase in sensitivity is desired. It must be said that recommendations

373  for using the filter DETAIL showed more inconsistencies compared to the other filters and this filter effect

374  could less clearly be linked to species traits.

375  We recommend using approved observations for species with a restricted range size, especially for large

376  birds. One noted exception was for the widespread species with a high classification error rate (profile 1),

377  where approved observations did impact model AUC positively.

378  When model AUC increased after filtering, sensitivity mostly increased and specificity decreased, with two

379  exceptions noted. First, species in profiles 1 and 2 were generally more difficult to identify, reflected by

380  either a high classification error rate (profile 1) or because they were small-bodied and unfamiliar to an

381  average observer (profile 2). For these profiles, we see that an increase in data quality by using either

382  filter could reduce the impact of false positives on model performance (i.e. increase specificity), except

383  for using approved observations for widespread species in profile 2. A side-effect was that sensitivity had

384     a greater potential to increase when sample size was reduced beyond the real data situation, even at high

385     reductions (Figure 3). A second exception, where specificity increased after filtering, was noted for familiar

386     species when using more detailed observations (profile 4) or data from active observers (profile 5). While

387     the positive impact of using only data from active observers on Δ specificity could be linked to higher

388     reporting probability (profile 5), the positive impact of using only detailed observations for familiar and

389     widespread species contradicted the negative association of Δ specificity with range size (Figure 2).

390

391 **4. Discussion**

392 In this study, we built recommendations for data quality filtering of opportunistic citizen science data

393 when used as input in species distribution models (SDMs), based on a set of a priori defined species traits.

394 Traits associated with a change in model performance after filtering were: body size, classification error

395 rate, familiarity, reporting probability and range size. Based on these traits, it was possible to generate

396 ecologically meaningful species profiles and make filtering recommendations (section 3.4). The analysis

397 of the species profiles mostly agreed with the results of a regression analysis but also gave new insights

398 on the relative importance of the different traits and trait combinations that lead to specific filtering

399 recommendations.

400 One of the main results was that, when choosing a quality filter, the taxonomic group a species belongs

401 to should be considered. This confirms previous findings based on the same dataset (Van Eupen et al.,

402 2021) and makes sense as taxonomic groups by default present differences in most of the considered

403 species traits due to differences in appearance, appeal, distribution etc. In an attempt to simplify the

404 results presented in this study, we have tested different approaches to generate the species profiles:

405 considering relative traits only, clustering of species for each taxonomic group separately and including

406 filter effects as active variables. Unfortunately, none of these approaches lead to profiles that were

407 ecologically more meaningful compared to the profiles suggested here (Table 2 and Table S1; evaluated

408 by species experts). Moreover, they lead to less consistent results (sections 2.3.1 and 2.3.2) or less explicit

409 filtering recommendations (i.e. larger confidence intervals in Figure 3). This confirms the expectation of

410 Van Eupen et al. (2021) who concluded that filtering recommendations can differ between taxonomic

411 groups, but that there might also be common traits among these groups that can refine them. The

412 selected approach indeed revealed that it is possible to formulate recommendations based on taxonomic

413 group and relative traits only (Table 2 and section 3.4). Absolute traits did not directly support

414 recommendations but aided the formation and interpretation of the species profiles as they either

415   characterized the most represented taxonomic group(s) or confirmed a profile's association with relative

416   traits.

417   The taxonomic bias towards bird species in citizen science data could explain some results, as it indicates

418   greater knowledge by the general public of this species group versus other groups such as butterflies and

419   dragonflies (Troudet et al., 2017; https://waarnemingen.be/stats/). As increased observer activity can

420   lead to higher experience and expertise (Johnston et al., 2017), this can explain why observer activity

421   mattered more for the less known taxonomic groups in this study (i.e. butterflies and dragonflies). For

422   example, experienced observers were better at detecting individuals of low-density insect populations

423   (Fitzpatrick et al., 2009) and increasing volunteer performance through training could reduce false positive

424   observations for pollinating insects (Ratnieks et al., 2016). These results can also be generalised to other

425   well-known taxonomic groups such as plants. Observer experience, for example, did not increase

426   volunteer performance for identifying an invasive plant species (Crall et al., 2011). Here, observers' self-

427   identified that comfort level was a better predictor of volunteer success.

428   The positive impact of using approved observations for birds, and especially for species with a restricted

429   range size, can be linked to the mechanism of record verification in the database (Swinnen et al., 2018),

430   whereby records that can be verified by photograph or sound play an important role. The verification

431   procedure consists of two main steps: (1) automated record validation by either image recognition or both

432   spatial and temporal proximity of new records to existing approved records and (2) manual expert

433   verification (when there is uncertainty in step 1). A decent photograph or sound record can thus easily

434   lead to multiple approved records and, by consequence, high photo or sound rates have more chance of

435   leading to approved (filtered) datasets of higher quality. Photo rates were, for example, generally higher

436   for bird and butterfly species with restricted range sizes (Table S1), which can explain why they benefitted

437   from using approved records. High photo or sound rates can also reduce the negative impact of locational

438   errors on model performance, especially for small sample sizes (Mitchell et al., 2017). Photographs are

439    often made from a closer distance, especially with the available easy-to-use identification apps (e.g.

440    ObsIdentify), leading to observations with lower locational uncertainty. When they are made from larger

441    distances, mostly for larger species (i.e. birds in this study), smartphone cameras will not suffice and an

442    observer needs a stronger camera lens. We believe that this is a pastime largely practised by more

443    experienced birders that are more likely to correctly register an individual's exact location compared to

444    an inexperienced observer. As for the importance of sound fragments, bird song usually indicates

445    territorial behaviour (Catchpole and Slater, 2008), hence observations made by sound are usually made

446    in birds' respective habitats. Additionally, the prevalence of locational errors in opportunistic bird data

447    will be larger compared to invertebrate species because of their high mobility (Maes et al., 2019), even at

448    a scale of 1 km², which was the resolution used in this study.

449    Large range size is associated with lower model performance because wide-ranged species usually occupy

450    a broad environmental niche and have less distinctive links with their habitat compared to species with a

451    restricted range size that usually have a narrow niche (e.g. Hernandez et al., 2006; Stockwell & Peterson,

452    2002). While increasing model performance for more widespread species through statistical methods or

453    survey design has been observed to be difficult (Brotons et al., 2007; Tessarolo et al., 2014), we observed

454    that using filtered data, especially from more active observers, had a positive impact on model

455    performance for widespread species. We argue, however, that range size in those cases is subordinate to

456    either classification error rate or the taxonomic group. Firstly, improving data quality is always important

457    for any species with a high misidentification risk (Table 2, profile 1). Misidentification errors can distort

458    estimates of species distributions (Costa et al., 2015; Cruickshank et al., 2019; Miller et al., 2011), even

459    though such errors were reduced by spatial aggregation of records (Kramer-Schadt et al., 2013; Van Eupen

460    et al., 2021). Misidentification risk has been found higher for species with similar physical appearance, for

461    example, because they are genetically related (Vantieghem et al., 2017) or have mimicking congeners

462    (Ratnieks et al., 2016). Secondly, widespread species in profile 4 are mostly large butterflies and, as

463 previously discussed, this taxonomic group might benefit more from using data from active observers.

464 Moreover, based on the relative traits only (i.e. not considering the taxonomic group) one would

465 intuitively assume that data quality filtering does not have such a pronounced positive impact in profile 4

466 because these widespread species were also more familiar and had lower error rates.

467 While retaining observations from active observers or approved observations showed clear associations

468 with taxonomic groups or relative species traits, retaining detailed observations showed more

469 inconsistencies, except for the positive impact on model performance for familiar species. Familiarity

470 might reflect the level of detail at which a species' ecology is known, hence data quality can be increased

471 by retaining more detailed observations for species that are familiar to an average citizen scientist.

472 Because retaining only detailed observations on average had the largest impact on sample size (Van Eupen

473 et al., 2021), the impact of sample size may be overruling the effect of the increase in data quality.

474 Reducing the sample size of presences generally impacts presence-only SDMs negatively as model

475 performance decreases, especially at low sample sizes, and performance variability increases (Hernandez

476 et al., 2006; Liu et al., 2019; van Proosdij et al., 2016). An increase in variability was mostly noted for

477 widespread species, as these species are more sensitive to small sample sizes (Liu et al., 2019). While large

478 reductions in sample size require attention, it remains important to realise that filtering simultaneously

479 increases data quality and thus model performance can also increase, especially when less than half of

480 the presences in a dataset are removed (Van Eupen et al., 2021).

481 Detectability did not appear to be an important trait in this study, while it has repeatedly been proven to

482 impact model performance positively (e.g. Pöyry et al., 2008; Seoane et al., 2005), and variation in

483 detectability is directly linked to the problem of imperfect detection in opportunistic presence-only data

484 (Dorazio, 2014). Species traits that are associated with increased detectability are, for example, high

485 abundance (Mccarthy et al., 2013), high singing rates (Sólymos et al., 2018), large body size (Johnston et

486 al., 2014; Pöyry et al., 2008), long lifespan and migratory behaviour (Carrascal et al., 2006). However, we

487    did not find proof that any of these traits were confounded with detectability in our analysis. One trait

488    that could have influenced the outcome for detectability was reporting probability because the way we

489    calculated reporting probability caused a moderate negative correlation between relative reporting

490    probability and relative detectability (Figure S1). However, reporting probability characterised only one

491    profile and thus implications for filtering recommendations would remain marginal.

492    While the highly fragmented (Antrop, 2004) and easily accessible landscape in our study region, Flanders,

493    has many benefits for studying species distributions, it was also one of the limitations. The largest benefit

494    was the consequent high spatial and temporal density of records in the *waarnemingen.be* database

495    (Herremans et al., 2018). On the other hand, because of the high density, the low importance of

496    detectability in our study could be an underestimation when studying regions with less fragmented and

497    larger conservation areas.

498    Another limitation was the insufficient availability of structured data for external model validation in the

499    original dataset (Van Eupen et al., 2021) leading to two restrictive features. First, data consisted of

500    relatively common species (minimum sample size was 432 presences). Rare habitat specialists from

501    habitats with restricted distribution ranges in Flanders (e.g. heathland) were thereby excluded from this

502    analysis. Since these are often targeted species in national and international biodiversity policy (De Ro et

503    al., 2021; Vanden Broeck et al., 2017), it would be useful to adjust the model validation strategy used in

504    Van Eupen et al. (2021) for those species to be able to formulate generic recommendations. Based on the

505    available data, building SDMs with validated data (for species with a restricted range size) or with data

506    from more active observers (for conspicuous invertebrates) could deliver the best results. Second, the

507    data showed sub-optimal representativeness of the taxonomic groups by the studied species. We argue,

508    however, that this imbalance in species representation is often inherent to opportunistic datasets (e.g.

509    over-representation of large birds in Callaghan et al., 2021).

510 Finally, some filter effects might have been impacted by the temporal and spatial aggregation of records

511 over the period 2014-2019 and in grid cells of 1x1 km. While a 1 km² resolution is a standard resolution in

512 Flemish biodiversity studies (e.g. Demolder et al., 2014; Rutten et al., 2019; Vantieghem et al., 2017),

513 performing the analysis at different scales might reveal higher or lower impacts of some traits.

514 **Conclusions**

515 Many have attempted to disentangle the relationships between species ecology and model performance,

516 and this study adds to that knowledge with some basic recommendations for data quality filtering for

517 three commonly studied taxonomic groups. Clustering species in species profiles based on traits that

518 resulted from a multiple regression analysis both highlighted the relative importance of species traits and

519 revealed new insights, and it is important to realise that one single trait does not necessarily predict a

520 species' response to filtering. We found that both the taxonomic group (more than absolute traits) and

521 relative species traits (rescaled values that can be compared among taxonomic groups) defined the impact

522 of data quality filtering on model performance. Our findings largely supported on: (1) the general species

523 knowledge among citizen scientists, with high importance of data quality for widespread and familiar

524 species in general and, more specifically, high importance of observer experience for less known

525 taxonomic groups; and (2) the mechanism of record verification in an opportunistic data platform, with

526 the high importance of submitting observations that can easily be verified, especially for species with

527 restricted range sizes. We encourage the further improvement of general species knowledge and

528 optimisation of record verification protocols in large citizen science projects. While adopting these

529 recommendations, it is always important to keep the goal of the study in mind (i.e. increasing model

530 discrimination capacity, sensitivity and/or specificity) and to keep an eye on the change in sample size

531 caused by stringent filtering.

532 **ACKNOWLEDGEMENTS**

533    We foremost thank the thousands of volunteers for collecting the millions of records that supported this

534    study; *Natuurpunt Studie* for making the data available for this research, in particular Joeri Belis and Karin

535    Gielen for composing the dataset; and Tim Adriaens for his help with the interpretation of the data. We

536    would also like to thank the anonymous reviewers for kindly improving the first version of this manuscript.

537    This work was supported by the Flemish Research Foundation FWO–SB [grant number 1S92118N].

538 **5. <u>References</u>**

539 Antrop, M., 2004. Landscape change and the urbanization process in Europe. Landsc. Urban Plan. 67, 9–
540     26. https://doi.org/10.1016/S0169-2046(03)00026-4

541 Bink, F.A., 1992. Ecologische atlas van de dagvlinders van Noordwest-Europa. Schuyt & Co Uitgevers en
542     Importeurs bv, Haarlem.

543 Bird, T.J., Bates, A.E., Lefcheck, J.S., Hill, N.A., Thomson, R.J., Edgar, G.J., Stuart-Smith, R.D.,
544     Wotherspoon, S., Krkosek, M., Stuart-Smith, J.F., Pecl, G.T., Barrett, N., Frusher, S., 2014. Statistical
545     solutions for error and bias in global citizen science datasets. Biol. Conserv. 173, 144–154.
546     https://doi.org/10.1016/j.biocon.2013.07.037

547 Brotons, L., Herrando, S., Pla, M., 2007. Updating bird species distribution at large spatial scales:
548     applications of habitat modelling to data from long-term monitoring programs. Divers. Distrib. 13,
549     276–288. https://doi.org/10.1111/j.1472-4642.2007.00339.x

550 Burgess, H.K., Debey, L.B., Froehlich, H.E., Schmidt, N., Theobald, E.J., Ettinger, A.K., Hillerislambers, J.,
551     Tewksbury, J., Parrish, J.K., 2017. The science of citizen science: Exploring barriers to use as a
552     primary research tool. Biol. Conserv. 208, 113–120. https://doi.org/10.1016/j.biocon.2016.05.014

553 Burnham, K.P., Anderson, D.R., Huyvaert, K.P., 2011. AIC model selection and multimodel inference in
554     behavioral ecology: some background, observations, and comparisons. Behav. Ecol. Sociobiol. 65,
555     23–35. https://doi.org/10.1007/s00265-010-1029-6

556 Callaghan, C.T., Poore, A.G.B., Hofmann, M., Roberts, C.J., Pereira, H.M., 2021. Large-bodied birds are
557     over-represented in unstructured citizen science data. Sci. Rep. 11, 19073.
558     https://doi.org/10.1038/S41598-021-98584-7

559 Carrascal, L.M., Javier, S., Palomino, D., Alonso, C.L., Lobo, J.M., 2006. Species-specific features affect the
560     ability of census-derived models to map winter avian distribution. Ecol. Res. 21, 681–691.
561     https://doi.org/10.1007/s11284-006-0173-y

562 Catchpole, C.K., Slater, P.J.B., 2008. Bird Song, 2nd ed. ed. Cambridge University Press, Cambridge.
563     https://doi.org/10.1017/CBO9780511754791

564 Chefaoui, R.M., Lobo, J.M., Hortal, J., 2011. Effects of species' traits and data characteristics on
565     distribution models of threatened invertebrates. Anim. Biodivers. Conserv. 34, 229–247.

566 Colwell, R.K., Xuan Mao, C., Chang, J., 2004. Interpolating, extrapolating, and comparing incidence-based
567     species accumulation curves. Ecology 85, 2717–2727.

568 Costa, H., Foody, G., Jiménez, S., Silva, L., 2015. Impacts of Species Misidentification on Species
569     Distribution Modeling with Presence-Only Data. ISPRS Int. J. Geo-Information 4, 2496–2518.
570     https://doi.org/10.3390/ijgi4042496

571 Crall, A.W., Newman, G.J., Stohlgren, T.J., Holfelder, K.A., Graham, J., Waller, D.M., 2011. Assessing
572     citizen science data quality: an invasive species case study. Conserv. Lett. 4, 433–442.
573     https://doi.org/10.1111/J.1755-263X.2011.00196.X

574 Cribari-Neto, F., Zeileis, A., 2010. Beta Regression in R. J. Stat. Softw. 34, 1–24.
575     https://doi.org/10.18637/jss.v034.i02

576  Cruickshank, S.S., Bühler, C., Benedikt, |, Schmidt, R., 2019. Quantifying data quality in a citizen science
577      monitoring program: False negatives, false positives and occupancy trends. Conserv. Sci. Pract. 1,
578      e54. https://doi.org/10.1111/CSP2.54

579  De Ro, A., Vanden broeck, A., Verschaeve, L., Van Dyck, H., Jacobs, I., T'Jollyn, F., Maes, D., 2021.
580      Occasional long distance dispersal does not prevent inbreeding in a threatened butterfly. BMC
581      Ecol. Evol. 21, 224.

582  Demolder, H., Schneiders, A., Spanhove, T., Maes, D., Van Landuyt, W., Adriaens, T., 2014. Hoofdstuk 4
583      Toestand biodiversiteit (INBO.R.2014.6194611), Natuurrapport - Toestand en trends van de
584      ecosystemen en ecosysteemdiensten in Vlaanderen. Instituut voor Natuur- en Bosonderzoek,
585      Brussels.

586  Dobson, A.D.M., Milner-Gulland, E.J., Aebischer, N.J., Beale, C.M., Brozovic, R., Coals, P., Critchlow, R.,
587      Dancer, A., Greve, M., Hinsley, A., Ibbett, H., Johnston, A., Kuiper, T., Le Comber, S., Mahood, S.P.,
588      Moore, J.F., Nilsen, E.B., Pocock, M.J.O., Quinn, A., Travers, H., Wilfred, P., Wright, J., Keane, A.,
589      2020. Making Messy Data Work for Conservation. One Earth 2, 455–465.
590      https://doi.org/10.1016/J.ONEEAR.2020.04.012

591  Dorazio, R.M., 2014. Accounting for imperfect detection and survey bias in statistical analysis of
592      presence-only data. Glob. Ecol. Biogeogr. 23, 1472–1484. https://doi.org/10.1111/GEB.12216

593  Ferrari, S., Cribari-Neto, F., Ferrari, S.L.P., 2004. Beta Regression for Modelling Rates and Proportions. J.
594      Appl. Stat. 31, 799–815. https://doi.org/10.1080/0266476042000214501

595  Fielding, A.H., Bell, J.F., 1997. A review of methods for the assessment of prediction errors in
596      conservation presence/absence models. Environ. Conserv. 24, 38–49.
597      https://doi.org/10.1017/S0376892997000088

598  Fitzpatrick, M.C., Preisser, E.L., Ellison, A.M., 2009. Observer Bias and the Detection of Low-Density
599      Populations. Ecol. Appl. 19, 1673–1679.

600  Guisan, A., Zimmermann, N.E., 2000. Predictive habitat distribution models in ecology. Ecol. Modell. 135,
601      147–186. https://doi.org/10.1016/S0304-3800(00)00354-9

602  Hanspach, J., Pompe, S., Klotz, S., 2010. Predictive performance of plant species distribution models
603      depends on species traits. Perspect. Plant Ecol. Evol. Syst. 12, 219–225.
604      https://doi.org/10.1016/j.ppees.2010.04.002

605  Henckel, L., Bradter, U., Jönsson, M., Isaac, N.J.B., Snäll, T., 2020. Assessing the usefulness of citizen
606      science data for habitat suitability modelling: Opportunistic reporting versus sampling based on a
607      systematic protocol. Divers. Distrib. 00, 1–15. https://doi.org/10.1111/ddi.13128

608  Hernandez, P.A., Graham, C.H., Master, L.L., Albert, D.L., 2006. The effect of sample size and species
609      characteristics on performance of different species distribution modeling methods. Ecography
610      (Cop.). 29, 773–785.

611  Herremans, M., Swinnen, K., Vanreusel, W., Vercayie, D., Veraghtert, W., Vanormelingen, P., 2018.
612      www.waarnemingen.be. Een veelzijdig portaal voor natuurgegevens. Natuur.focus 17, 153–166.

613  Husson, F., Josse, J., Pagès, J., 2010. Principal component methods - hierarchical clustering - partitional
614      clustering: why would we need to choose for visualizing data? Applied Mathematics Department.

615    Isaac, N.J.B., Jarzyna, M.A., Keil, P., Dambly, L.I., Boersch-Supan, P.H., Browning, E., Freeman, S.N.,
616        Golding, N., Guillera-Arroita, G., Henrys, P.A., Jarvis, S., Lahoz-Monfort, J., Pagel, J., Pescott, O.L.,
617        Schmucki, R., Simmonds, E.G., O'Hara, R.B., 2020. Data Integration for Large-Scale Models of
618        Species Distributions. Trends Ecol. Evol. 35, 56–67. https://doi.org/10.1016/j.tree.2019.08.006

619    Isaac, N.J.B., Pocock, M.J.O., 2015. Bias and information in biological records. Biol. J. Linn. Soc. 115, 522–
620        531. https://doi.org/10.1111/bij.12532

621    Jiménez-Valverde, A., 2012. Insights into the area under the receiver operating characteristic curve
622        (AUC) as a discrimination measure in species distribution modelling. Glob. Ecol. Biogeogr. 21, 498–
623        507. https://doi.org/10.1111/j.1466-8238.2011.00683.x

624    Johnston, A., Fink, D., Hochachka, W.M., Kelling, S., 2017. Estimates of observer expertise improve
625        species distributions from citizen science data. Methods Ecol. Evol. 00, 1–10.
626        https://doi.org/10.1111/2041-210X.12838

627    Johnston, A., Hochachka, W.M., Strimas-Mackey, M.E., Gutierrez, V.R., Robinson, O.J., Miller, E.T., Auer,
628        T., Kelling, S.T., Fink, D., 2021. Analytical guidelines to increase the value of community science
629        data: An example using eBird data to estimate species distributions. Divers. Distrib. 27, 1265–1277.
630        https://doi.org/10.1111/DDI.13271

631    Johnston, A., Newson, S.E., Risely, K., Musgrove, A.J., Massimino, D., Baillie, S.R., Pearce-Higgins, J.W.,
632        2014. Species traits explain variation in detectability of UK birds. Bird Study 61, 340–350.
633        https://doi.org/10.1080/00063657.2014.941787

634    Kamp, J., Oppel, S., Heldbjerg, H., Nyegaard, T., Donald, P.F., 2016. Unstructured citizen science data fail
635        to detect long-term population declines of common birds in Denmark. Divers. Distrib. 22, 1024–
636        1035. https://doi.org/10.1111/ddi.12463

637    Kosmala, M., Wiggins, A., Swanson, A., Simmons, B., 2016. Assessing data quality in citizen science.
638        Front. Ecol. Environ. 14, 551–560. https://doi.org/10.1002/fee.1436

639    Kramer-Schadt, S., Niedballa, J., Pilgrim, J.D., Schröder, B., Lindenborn, J., Reinfelder, V., Stillfried, M.,
640        Heckmann, I., Scharf, A.K., Augeri, D.M., Cheyne, S.M., Hearn, A.J., Ross, J., Macdonald, D.W.,
641        Mathai, J., Eaton, J., Marshall, A.J., Semiadi, G., Rustam, R., Bernard, H., Alfred, R., Samejima, H.,
642        Duckworth, J.W., Breitenmoser-Wuersten, C., Belant, J.L., Hofer, H., Wilting, A., 2013. The
643        importance of correcting for sampling bias in MaxEnt species distribution models. Divers. Distrib.
644        19, 1366–1379. https://doi.org/10.1111/DDI.12096

645    Le, S., Josse, J., Husson, F., 2008. FactoMineR: An R Package for Multivariate Analysis. J. Stat. Softw. 25,
646        1–18.

647    Lindenmayer, D., Woinarski, J., Legge, S., Southwell, D., Lavery, T., Robinson, N., Scheele, B., Wintle, B.,
648        2020. A checklist of attributes for effective monitoring of threatened species and threatened
649        ecosystems. J. Environ. Manage. 262, 110312. https://doi.org/10.1016/j.jenvman.2020.110312

650    Liu, C., Newell, G., White, M., 2019. The effect of sample size on the accuracy of species distribution
651        models: considering both presences and pseudo-absences or background sites. Ecography (Cop.).
652        42, 535–548. https://doi.org/10.1111/ecog.03188

653    Lobo, J.M., Jiménez-valverde, A., Real, R., 2008. AUC: A misleading measure of the performance of
654        predictive distribution models. Glob. Ecol. Biogeogr. 17, 145–151. https://doi.org/10.1111/j.1466-
655        8238.2007.00358.x

656    MacKenzie, D.I., Nichols, J.D., Royle, J.A., Pollock, K.H., Bailey, L., Hines, J.E., 2017. Occupancy Estimation
657           and Modeling: Inferring Patterns and Dynamics of Species Occurrence. Elsevier.

658    Maes, D., Bauwens, D., De Bruyn, L., Anselin, A., Vermeersch, G., Van Landuyt, W., De Knijf, G., Gilbert,
659           M., 2005. Species richness coincidence: Conservation strategies based on predictive modelling.
660           Biodivers. Conserv. 14, 1345–1364. https://doi.org/10.1007/S10531-004-9662-X

661    Maes, D., Brosens, D., Desmet, P., Piesschaert, F., Van Hoey, S., Adriaens, T., Dekoninck, W., Devos, K.,
662           Lock, K., Onkelinx, T., Packet, J., 2019. A database of threat statuses and life-history traits of Red
663           List species in Flanders (northern Belgium). Biodivers. Data J. 7, e34089.
664           https://doi.org/10.3897/BDJ.7.e34089

665    Maes, D., Isaac, N.J.B., Harrower, C.A., Collen, B., van Strien, A.J., Roy, D.B., 2015. The use of
666           opportunistic data for IUCN Red List assessments. Biol. J. Linn. Soc. 115, 690–706.
667           https://doi.org/10.1111/bij.12530

668    Matutini, F., Baudry, J., Pain, G., Sineau, M., Pithon, J., 2021. How citizen science could improve species
669           distribution models and their independent assessment. Ecol. Evol. 11, 3028–3039.
670           https://doi.org/10.1002/ece3.7210

671    McCarthy, M.A., Moore, J.L., Morris, W.K., Parris, K.M., Garrard, G.E., Vesk, P.A., Rumpff, L., Giljohann,
672           K.M., Camac, J.S., Bau, S.S., Friend, T., Harrison, B., Yue, B., 2013. The influence of abundance on
673           detectability. Oikos 122, 717–726. https://doi.org/10.1111/j.1600-0706.2012.20781.x

674    McPherson, J.M., Jetz, W., Rogers, D.J., 2004. The effects of species' range sizes on the accuracy of
675           distribution models: Ecological phenomenon or statistical artefact? J. Appl. Ecol. 41, 811–823.
676           https://doi.org/10.1111/J.0021-8901.2004.00943.X

677    Menard, S., 2001. Applied Logistic Regression Analysis. 2nd edition. SAGE Publications, Inc.

678    Miller, D.A., Nichols, J.D., McClintock, B.T., Grant, E.H.C., Bailey, L.L., Weir, L.A., 2011. Improving
679           occupancy estimation when two types of observational error occur: non-detection and species
680           misidentification. Ecology 92, 1422–1428. https://doi.org/10.1890/10-1396.1

681    Mitchell, P.J., Monk, J., Laurenson, L., 2017. Sensitivity of fine-scale species distribution models to
682           locational uncertainty in occurrence data across multiple sample sizes. Methods Ecol. Evol. 8, 12–
683           21. https://doi.org/10.1111/2041-210X.12645

684    Morton, E.S., 1975. Ecological Sources of Selection on Avian Sounds. Am. Nat. 109, 17–34.

685    Newbold, T., Hudson, L.N., Hill, S.L.L., Contu, S., Lysenko, I., Senior, R.A., Börger, L., Bennett, D.J.,
686           Choimes, A., Collen, B., Day, J., De Palma, A., Díaz, S., Echeverria-Londoño, S., Edgar, M.J., Feldman,
687           A., Garon, M., Harrison, M.L.K., Alhusseini, T., Ingram, D.J., Itescu, Y., Kattge, J., Kemp, V.,
688           Kirkpatrick, L., Kleyer, M., Laginha, D., Correia, P., Martin, C.D., Meiri, S., Novosolov, M., Pan, Y.,
689           Phillips, H.R.P., Purves, D.W., Robinson, A., Simpson, J., Tuck, S.L., Weiher, E., White, H.J., Ewers,
690           R.M., Mace, G.M., Scharlemann, J.P.W., Purvis, A., 2015. Global effects of land use on local
691           terrestrial biodiversity. Nature 520, 45–50. https://doi.org/10.1038/nature14324

692    Phillips, S.J., Anderson, R.P., Schapire, R.E., 2006. Maximum entropy modeling of species geographic
693           distributions. Ecol. Modell. 190, 231–259. https://doi.org/10.1016/j.ecolmodel.2005.03.026

694    Pöyry, J., Luoto, M., Heikkinen, R.K., Saarinen, K., 2008. Species traits are associated with the quality of
695           bioclimatic models. Glob. Ecol. Biogeogr. 17, 403–414. https://doi.org/10.1111/j.1466-

696        8238.2007.00373.x

697    R Core Team, 2021. R: A language and environment for statistical computing.

698    Ratnieks, F.L.W., Schrell, F., Sheppard, R.C., Brown, E., Bristow, O.E., Garbuzov, M., 2016. Data reliability
699        in citizen science: learning curve and the effects of training method, volunteer background and
700        experience on identification accuracy of insects visiting ivy flowers. Methods Ecol. Evol. 7, 1226–
701        1235. https://doi.org/10.1111/2041-210X.12581

702    Rutten, A., Casaer, J., Swinnen, K.R.R., Herremans, M., Leirs, H., 2019. Future distribution of wild boar in
703        a highly anthropogenic landscape: Models combining hunting bag and citizen science data. Ecol.
704        Modell. 411, 108804. https://doi.org/10.1016/j.ecolmodel.2019.108804

705    Seoane, J., Carrascal, L.M., Alonso, L., Palomino, D., 2005. Species-specific traits associated to prediction
706        errors in bird habitat suitability modelling. Ecol. Modell. 185, 299–308.
707        https://doi.org/10.1016/j.ecolmodel.2004.12.012

708    Serra-Diaz, J.M., Enquist, B.J., Maitner, B., Merow, C., Svenning, J.C., 2017. Big data of tree species
709        distributions: how big and how good? For. Ecosyst. 4. https://doi.org/10.1186/s40663-017-0120-0

710    Sólymos, P., Matsuoka, S.M., Stralberg, D., Barker, N.K.S., Bayne, E.M., 2018. Phylogeny and species
711        traits predict bird detectability. Ecography (Cop.). 41, 1595–1603.
712        https://doi.org/10.1111/ecog.03415

713    Steen, V.A., Elphick, C.S., Tingley, M.W., 2019. An evaluation of stringent filtering to improve species
714        distribution models from citizen science data. Biodivers. Res. 25, 1857–1869.
715        https://doi.org/10.1111/ddi.12985

716    Stockwell, D.R.B., Peterson, A.T., 2002. Effects of sample size on accuracy of species distribution models.
717        Ecol. Modell. 148, 1–13.

718    Storchová, L., Hořák, D., 2018. Life-history characteristics of European birds. Glob. Ecol. Biogeogr. 27,
719        400–406. https://doi.org/10.1111/geb.12709

720    Sullivan, B.L., Wood, C.L., Iliff, M.J., Bonney, R.E., Fink, D., Kelling, S., 2009. eBird: A citizen-based bird
721        observation network in the biological sciences. Biol. Conserv. 142, 2282–2292.
722        https://doi.org/10.1016/j.biocon.2009.05.006

723    Swinnen, K.R.R., Jacobs, A., Claus, K., Ruyts, S., Vercayie, D., Lambrechts, J., Herremans, M., n.d. 'Animals
724        under wheels': wildlife roadkill data collection by citizen scientists as a part of their nature
725        recording activities. Nat. Conserv.

726    Swinnen, K.R.R., Vercayie, D., Vanreusel, W., Barendse, R., Boers, K., Bogaert, J., Dekeukeleire, D.,
727        Driessens, G., Dupriez, P., Jooris, R., Steeman, R., van Asten, K., van den Neucker, T., van
728        Dorsselaer, P., van Vooren, P., Wysmantel, N., Gielen, K., Desmet, P., Herremans, M., 2018.
729        Waarnemingen.be-Non-native plant and animal occurrences in Flanders and the Brussels Capital
730        Region, Belgium. BioInvasions Rec. 7, 335–342. https://doi.org/10.3391/bir.2018.7.3.17

731    Tessarolo, G., Rangel, T.F., Araújo, M.B., Hortal, J., 2014. Uncertainty associated with survey design in
732        Species Distribution Models. Divers. Distrib. 20, 1258–1269. https://doi.org/10.1111/ddi.12236

733    Thomaes, A., Kervyn, T., Maes, D., 2008. Applying species distribution modelling for the conservation of
734        the threatened saproxylic Stag Beetle (Lucanus cervus). Biol. Conserv. 141, 1400–1410.

735        https://doi.org/10.1016/j.biocon.2008.03.018

736    Troudet, J., Grandcolas, P., Blin, A., Vignes-Lebbe, R., Legendre, F., 2017. Taxonomic bias in biodiversity
737        data and societal preferences. Sci. Rep. 7, 9132. https://doi.org/10.1038/s41598-017-09084-6

738    Urban, M.C., Bocedi, G., Hendry, A.P., Mihoub, J.B., Pe'er, G., Singer, A., Bridle, J.R., Crozier, L.G., De
739        Meester, L., Godsoe, W., Gonzalez, A., Hellmann, J.J., Holt, R.D., Huth, A., Johst, K., Krug, C.B.,
740        Leadley, P.W., Palmer, S.C.F., Pantel, J.H., Schmitz, A., Zollner, P.A., Travis, J.M.J., 2016. Improving
741        the forecast for biodiversity under climate change. Science (80-. ). 353, aad8466.
742        https://doi.org/10.1126/science.aad8466

743    Van Eupen, C., Maes, D., Herremans, M., Swinnen, K.R.R., Somers, B., Luca, S., 2021. The impact of data
744        quality filtering of opportunistic citizen science data on species distribution model performance.
745        Ecol. Modell. 444, 109453. https://doi.org/10.1016/j.ecolmodel.2021.109453

746    van Proosdij, A.S.J., Sosef, M.S.M., Wieringa, J.J., Raes, N., 2016. Minimum required number of specimen
747        records to develop accurate species distribution models. Ecography (Cop.). 39, 542–552.
748        https://doi.org/10.1111/ECOG.01509

749    Van Strien, A.J., Van Swaay, C.A.M., Termaat, T., 2013. Opportunistic citizen science data of animal
750        species produce reliable estimates of distribution trends if analysed with occupancy models. J.
751        Appl. Ecol. 50, 1450–1458. https://doi.org/10.1111/1365-2664.12158

752    Vanden Broeck, A., Maes, D., Kelager, A., Wynhoff, I., WallisDeVries, M.F., Nash, D.R., Oostermeijer,
753        J.G.B., Van Dyck, H., Mergeay, J., 2017. Gene flow and effective population sizes of the butterfly
754        Maculinea alcon in a highly fragmented, anthropogenic landscape. Biol. Conserv. 209, 89–97.
755        https://doi.org/10.1016/j.biocon.2017.02.001

756    Vantieghem, P., Maes, D., Kaiser, A., Merckx, T., 2017. Quality of citizen science data and its
757        consequences for the conservation of skipper butterflies (Hesperiidae) in Flanders (northern
758        Belgium). J. Insect Conserv. 21, 451–463. https://doi.org/10.1007/s10841-016-9924-4

759    Vermeersch, G., Devos, K., Driessens, G., Evereaert, J., Feys, S., Herremans, M., Onkelinx, T., Stienen,
760        E.W.M., T'Jollyn, F., Anselin, A., 2020. Broedvogels in Vlaanderen 2013-2018. Medelingen van het
761        Instituut voor Natuur- en Bosonderzoek 2020 (1), Brussel.
762        https://doi.org/10.21436/inbor.18794135

763    Żmihorski, M., Dziarska-Pałac, J., Sparks, T.H., Tryjanowski, P., 2013. Ecological correlates of the
764        popularity of birds and butterflies in Internet information resources. Oikos 122, 183–190.
765        https://doi.org/10.1111/J.1600-0706.2012.20486.X

766    Zurell, D., Franklin, J., König, C., Bouchet, P.J., Dormann, C.F., Elith, J., Fandos, G., Feng, X., Guillera-
767        Arroita, G., Guisan, A., Lahoz-Monfort, J.J., Leitão, P.J., Park, D.S., Peterson, A.T., Rapacciuolo, G.,
768        Schmatz, D.R., Schröder, B., Serra-Diaz, J.M., Thuiller, W., Yates, K.L., Zimmermann, N.E., Merow,
769        C., 2020. A standard protocol for reporting species distribution models. Ecography (Cop.). 43,
770        1261–1277. https://doi.org/10.1111/ECOG.04960

771     [dataset] Van Eupen, C., Maes, D., Herremans, M., Swinnen, K.R.R., Somers, B., Luca, S., 2021b. The
772        impact of data quality filtering of opportunistic citizen science data on species distribution model
773        performance: dataset used for Maxent modelling., Dryad, Dataset,
774        https://doi.org/10.5061/dryad.jwstqjq83

775