1 # Treating gaps and biases in biodiversity data as a missing data
2 # problem

3

4 Diana E. Bowler[1]*, Robin J. Boyd[1], Corey T. Callaghan[2], Robert A. Robinson[3], Nick J. B.

5 Isaac[1], Michael J. O. Pocock[1]

6

7 Affiliations

8 1 UK Centre for Ecology &amp; Hydrology, Wallingford, OX10 8BB, UK

9 2 Department of Wildlife Ecology and Conservation, Fort Lauderdale Research and Education

10 Center, University of Florida, Davie, FL 33314-7719

11 3 British Trust for Ornithology, The Nunnery, Thetford, IP24 2PU, UK

12

13 * Corresponding author: diana.e.bowler@gmail.com

14

# Abstract

Big biodiversity datasets have great potential for monitoring and research because of their large taxonomic, geographic and temporal scope. Such datasets have become especially important for assessing the temporal change of species' populations and distributions. Gaps in the available data, however, often hinder drawing large-scale inferences about species' trends. Here, we conceptualise biodiversity data gaps as a missing data problem, which provides a unifying framework for the challenges and potential solutions across different types of biodiversity datasets. We characterise the typical types of data gaps in biodiversity data as different classes of missing data and then use missing data theory to explore the implications for different research questions. By using this framework, we show that bias due to data gaps can arise when the factors affecting sampling and/or data availability overlap with those affecting biodiversity. But the outcome also depends on the ecological questions, which determines choices around the analytical approach. We argue that typical approaches to long-term species trend modelling are especially susceptible to data gaps since such models do not tend to account for the factors that drive missingness. To identify general solutions, we review empirical studies and use simulation studies to compare some of the most frequently employed approaches to deal with data gaps, including subsampling, weighting and imputation. All these methods have the potential to reduce bias but may come at the cost of increased uncertainty of parameter estimates. Weighting approaches are arguably the least used so far in ecology and have the potential to reduce both the bias and variance of parameter estimates. Regardless of the method, the ability to reduce bias critically depends on knowledge of, and the availability of data on, the factors creating data gaps. We use our review to outline the necessary considerations when dealing with data gaps at different stages of the data collection and analysis workflow.

38    Keywords: Biodiversity change; Citizen Science; Ecological Modelling; Macroecology; Spatial

39    bias

40

Contents

# I. Introduction: uneven sampling of biodiversity

Ecologists have ever-growing access to data on species' occurrence and abundances. Potential sources of data include long-term citizen-science monitoring schemes (such as the North American Breeding Bird Survey) (Bled *et al.*, 2013), data aggregators (such as the Global Biodiversity Information Facility) (Garcia-Rosello *et al.*, 2015), remote-sensing platforms (Fretwell, Scofield & Phillips, 2017) and synthesis databases (such as BioTIME or the Living Planet Database) (Dornelas *et al.*, 2014). Since these data cover broad spatial and temporal scales, they are especially useful for large-scale questions, for instance, about species' distributions, population and community-level trends, and ecological niches (Chandler *et al.*, 2017; Sullivan *et al.*, 2017; Fink *et al.*, 2020). These data also underpin many biodiversity trend indicators that are central for national and international conservation policy (Gregory *et al.*, 2005; van Swaay *et al.*, 2008; Fraisl *et al.*, 2020).

Despite the impressive volume of data, biodiversity data, regardless of the source, tend to be filled with gaps and redundancies (Boakes *et al.*, 2010). Data gaps are not necessarily problematic; indeed, most studies rely on statistical inference to make inferences about a broader region of interest from a sample. Data gaps, however, can be problematic when they lead to biases (Boakes *et al.*, 2010; Bled *et al.*, 2013; Amano, Lamming & Sutherland, 2016). Already many ecologists have raised concerns about the impacts of bias on estimated spatial or temporal biodiversity patterns (Bayraktarov *et al.*, 2019; Valdez *et al.*, 2023). Developing methods to deal with data gaps and associated biases within large-scale biodiversity data is an increasingly important task to make full use of the growing big data sources.

87     Patterns in the availability of biodiversity data can be affected by the original motivations

88     for, and constraints on, data collection activities. While some data are collected as part of

89     scientific studies, much of the available data on species' occurrence and abundance are collected

90     through citizen science initiatives (Chandler *et al.*, 2017). Spatial patterns in data availability

91     from citizen science have been especially well-studied. Citizen science programs have varying

92     degrees of protocol and sampling designs (Isaac & Pocock, 2015; Pocock *et al.*, 2017) but more

93     data are typically collected in accessible areas such as near roads and urban areas (Geldmann *et*

94     *al.*, 2016). Such biases are not unique to citizen science data, as even data collected during

95     formal scientific studies have potential sampling biases towards regions undergoing less habitat

96     change, which may lead to underestimates of biodiversity change (Gonzalez *et al.*, 2016; Forister

97     *et al.*, 2023; Cardinale *et al.*, 2018). Various solutions have been proposed to deal with these

98     biases (Hefley *et al.*, 2013; Cretois *et al.*, 2021; Johnston *et al.*, 2020; Ver Hoef *et al.*, 2021), but

99     there is still a lack of a general framework for ecologists to guide decisions on when and how to

100    deal with data gaps.

101     Here, we show how using missing data theory (Rubin, 1976) can unify problems

102    associated with data gaps across different types of biodiversity datasets. Missing data are a

103    widespread problem crossing disciplines, with a large body of literature on the implications and

104    possible solutions (Little & Rubin, 2019; Carpenter & Kenward, 2012). We expect that aligning

105    the generalized problem of missing data, conceptualized within missing data theory, to the

106    problem of biodiversity data biases discussed above will yield opportunities so far overlooked.

107    We mostly focus our review on modelling trends in species occupancy or abundance using

108    monitoring data collected by volunteer citizen scientists, but the ideas transfer to other types of

109    biodiversity data or questions. We show that bias is not a property of a dataset but rather a

110    property of the use of a dataset for a specific question and target population that are imposed by

111    the data analyst. We review some commonly used solutions to missing data to highlight potential

112    approaches that could be considered in biodiversity analyses.

## 113    II. Classifying data gaps using missing data theory

### 114    (1) Biodiversity data gaps

115    Species occurrence or abundance data can have gaps in different dimensions. We distinguish

116    between spatial, annual and within-year gaps (Fig. 1). We define spatial gaps as those formed by

117    sites with no data, and annual gaps as those formed by a lack of data in some years at sites that

118    have been otherwise sampled. Together, spatial and annual gaps determine the spatial and

119    temporal coverage of a dataset. Within-year gaps arise when data are lacking in specific seasons

120    or months, which can be important because most organisms are seasonal and multiple visits are

121    usually necessary to robustly estimate detection probabilities. Considering why these gaps arise

122    can help understand their likely impact, for instance, on long-term species trend estimation. Data

123    gaps are found in different types of monitoring data including highly structured monitoring

124    schemes with a standardised protocol, such as many national bird survey schemes, as well as

125    opportunistic monitoring data that are typically an aggregation of heterogeneous observations.

126    Biodiversity datasets can also have taxonomic gaps (Troudet *et al.*, 2017), but this is outside the

127    scope of this paper.

128

**Fig. 1 Different types of data gaps within biodiversity data.**

We imagine a scenario where there are multiple survey visits across sites and years. Visits can be in

response to a protocol ('structured' data) or opportunistic ('unstructured'), and repeat visits can be by the

same or multiple recorders. Data gaps, or more generally uneven data availability, can arise due to (a)

within-year gaps (e.g., blue square, i.e., ordinarily there are three visits, but some sites are only visited

once or twice in a year); (b) annual gaps (e.g., yellow square, i.e., some sites that are usually sampled are

entirely unvisited in some years) or (c) spatial gaps (e.g., red square, i.e., some sites within the region of

interest are never visited across all years). Some sites are well-sampled within and across years and hence

have no missing data (e.g., green square).

138

While both structured and opportunistic monitoring data can be affected by similar data

gaps (Binley & Bennett, 2023), there are some key differences between these types of

monitoring data. Moreover, structured schemes themselves vary in the degree of structure and

standardisation. In structured schemes with a formal spatial sampling design, data gaps include

7

143　　both planned and unplanned gaps. Planned gaps arise because only a sample of sites was ever

144　　intended to be sampled. Unplanned gaps occur because of failure to recruit and retain surveyors

145　　at sites that were intended to be sampled (Zhang *et al.*, 2021; Marsh & Cosentino, 2019). In most

146　　other types of data, gaps are neither planned or unplanned. Some monitoring schemes have

147　　sampling protocols but participants are free to choose their own sampling sites. In fully

148　　opportunistic monitoring schemes, participants make individual decisions about where to sample

149　　and gaps emerge from unevenness in the cumulative sampling effort of all participants. Due to

150　　the high number of participants, and lack of coordination of their effects, sampling effort is

151　　generally more strongly skewed across space and time in opportunistic schemes than in

152　　structured schemes, leading to more pervasive data gaps (Geldmann *et al.*, 2016). Synthesis

153　　databases such as BioTIME and the Living Planet Database, and data aggregators such as GBIF,

154　　are similar in these respects to schemes without a formal spatial sampling design since they

155　　contain data that were independently collected as part of separate studies, without coordinated

156　　efforts.

157　　　　Despite these differences, correlates of data gaps tend to be similar across monitoring

158　　schemes, especially those involving citizen scientists. Spatial gaps often occur in remote areas

159　　because there is a smaller pool of potential participants nearby (Geldmann *et al.*, 2016;

160　　Mandeville, Nilsen & Finstad, 2022). Spatial gaps can also be more common where species have

161　　lower abundance or land cover is perceived to be less attractive for biodiversity and for

162　　surveying e.g., agricultural land (Tulloch *et al.*, 2013; Dambly *et al.*, 2021; Marsh & Cosentino,

163　　2019). Annual gaps can arise due to project turnover or because of external factors (e.g. the 2020

164　　season for most countries was highly compromised by the Covid-19 pandemic). Annual gaps

165　　have also been linked with local land use changes that negatively affected species abundance and

166  the attractiveness of a site for sampling (Zhang *et al.*, 2021; Marsh & Cosentino, 2019). Within-

167  year data gaps can be caused by periods of inclement weather (Zimney & Smart, 2022; Diekert

168  *et al.*, 2023) or vary seasonally e.g., missing surveys for butterflies are more common at start and

169  end of the main flight period (Dennis *et al.*, 2016), while bird sampling can be higher during

170  their migration periods (La Sorte & Somveille, 2020).

171

## (2) Classes of missing data

173  Within the classic missing data theory, there are three classes of missing data (Missing

174  Completely at Random, Missing at Random, Missing Not at Random), defined below, each with

175  different consequences for bias (Table 1) (Rubin, 1976; Nakagawa & Freckleton, 2008; Little &

176  Rubin, 2019). These classes vary in their missing data mechanism, which describes the

177  relationship between the probability of missing data (or sampling effort in the monitoring

178  context) and the values of other variables. Hefley et al. (2013) already proposed viewing spatial

179  biases in presence-only data as a form of missing data. Here, we extend it more broadly across

180  different types of biodiversity data.

181      Within the context of biodiversity data, missingness can be regarded as Missing

182  Completely at Random (MCAR) if the factors affecting biodiversity sampling, and causing

183  missingness, are independent of those affecting biodiversity (Table 1). Under MCAR, the

184  observed data are effectively a random sample of the whole population, and the values of the

185  variable of interest are similar in sampled and non-sampled sites or times. For instance, if site

186  selection is driven by human accessibility, but species distribution is primarily driven by climate,

187  and if accessibility and climate are not correlated, then spatial data gaps would be MCAR.

188  Within-year gaps associated with weekdays (Evans & Day, 2002; Courter *et al.*, 2013), or annual

9

189  gaps associated with project turnover, are also examples likely to cause MCAR data gaps since

190  such gaps are probably not associated with biodiversity patterns (Table 1). In this case, missing

191  data could reduce the precision of parameters estimates through reduced sample size, but not

192  increase the bias.

193      When the factors affecting sampling are the same as, or correlated with, those affecting

194  biodiversity, the missing data mechanism can either be Missing at Random (MAR) or Missing

195  Not at Random (MNAR). For instance, if road density affects both sampling probability and

196  species abundance, then spatial gaps are not MCAR. Road density might affect sampling

197  probability directly (e.g., if people are more often looking for wildlife along roads) or indirectly

198  (e.g., if road density affects species detectability); in either case, road density influences data

199  gaps. Similarly, habitat degradation could reduce both species abundance and observer retention

200  to continue sampling at a site, creating an annual data gap that is MAR or MNAR (Table 1). In

201  these cases, there are systematic differences in the biodiversity quantity of interest between

202  sampled and non-sampled sites or times (Table 1).

203      To borrow from an infamous quote, if we regard data gaps as "unknowns", then MAR

204  can be thought of as "known unknowns" while MNAR are "unknown unknowns".  The "known"

205  needed for MAR is knowledge and availability of data on the shared covariates affecting

206  sampling and biodiversity.  If complete data for shared covariates are available and included in

207  the analysis, then the missing data mechanism is MAR. Hence, despite its name, MAR does not

208  mean that sampling effort is randomly distributed in the landscape. Rather, it means that the

209  covariates affecting sampling are known and that there is available covariate data to fully explain

210  the differences between sampled and non-sampled potential data. If any of the relevant factors

211  affecting sampling and biodiversity are unknown, or not modelled, the missing data mechanism

212    becomes MNAR (Table 1). Hence, decisions of the analyst can determine whether a data gap is

213    MNAR or MAR (discussed more fully in section III). MNAR may also arise when missingness

214    is dependent on the value of biodiversity itself, i.e., if sampling effort directly depends on species

215    occurrence or abundance.


216         Statistical tests can only partly indicate which missing data class is most likely (Little,

217    1988). Analysis of relationships between data availability and observed covariates can point

218    towards MAR if some relationships are significant. But a lack of any association, or an

219    incomplete explanation of data gaps, could reflect MCAR or MNAR. Because MNAR is

220    associated with unavailable data, it cannot be directly tested. Concerns about whether

221    missingness in the biodiversity data is directly associated with its values could be explored if

222    there is a related variable that is fully available (Wu, 2022). We argue that MCAR is unlikely in

223    most biodiversity data since unplanned data gaps can affect even the most structured monitoring

224    schemes.


225

226 **Table 1 Missing data mechanisms in biodiversity data, including examples and implications**

| Mechanism | Typical meaning | Meaning in the context of biodiversity data | Examples | Typical implications |
|---|---|---|---|---|
| Missing completely at random (MCAR) | Missingness is independent of observed and unobserved variables. | Sampling is independent of any covariates, or covariates that affect sampling probability are independent of those affecting biodiversity | Within-year: Weekday gaps<br>Annual/Spatial: Gaps caused by the completion of a fixed-term project or retirement of a participant | Ignorable |
| Missing at random (MAR) | Missingness is associated with observed data but not any unobserved variables | Covariates that affect sampling probability are shared with those affecting biodiversity, but data are available on all these covariates | within-year: Season (day of year)<br>Annual: Urban development<br>Spatial: Accessibility | Ignorable if the model includes all relevant covariates |
| Missing not at random (MNAR) | Missingness depends on unobserved variables or the missing values itself | CS sampling varies with biodiversity value or an unknown or unavailable covariate affects sampling and biodiversity | within-year/annual/spatial: unknown factors causing variation in species activity/abundance that are also correlated with sampling effort | Non-ignorable - the missing data mechanism needs to be modelled |

# III. Implications of missingness for ecological questions

228 Missing data (i.e., data gaps) themselves do not necessarily have strong impacts on the results of

229 biodiversity modelling, but can depend on the specific question and parameter of interest

230 (Bartlett, Harel & Carpenter, 2015; Collins, Schafer & Kam, 2001; Little *et al.*, 2022). Viewing

231 data gaps as a form of missing data can help decide whether a particular data gap matters. As we

232  note above, data gaps that are MCAR do not cause bias, but data gaps in biodiversity data are

233  unlikely to be wholly MCAR. For a data gap to be MAR rather than MNAR can depend on the

234  ecological question being pursued by an analyst. This is because the 'missing at random'

235  assumption of MAR is conditional on controlling for covariates affecting sampling probability,

236  which means that these covariates are known, reflected in available data and included in the

237  analysis (Fig. 2) (Conn, Thorson & Johnson, 2017; Hefley *et al.*, 2013). Different ecological

238  questions will lead to different decisions about which variables to include in an analysis. Hence,

239  data gaps of the same dataset might be MAR under some questions but MNAR under others. To

240  illustrate these potential differences, we contrast two typical questions asked with biodiversity

241  data.

242

## (1) Understanding the roles of environmental drivers on species' distributions

244  Monitoring data are often used to understand the environmental factors explaining species

245  distribution patterns. The implications of missing data for species distribution models have been

246  considered in terms of niche truncation. Niche truncation happens when a dataset only contains

247  occurrence data from part of the geographic range of a species, which usually also means that the

248  dataset only covers part of the ecological/environmental space that is suitable for the species

249  (Chevalier *et al.*, 2022; Albert *et al.*, 2010; Guo *et al.*, 2023). These studies show that the

250  implications of niche truncation depend on the functional form of the relationship between the

251  associated covariate and the species response (Chevalier *et al.*, 2022) and whether occurrence

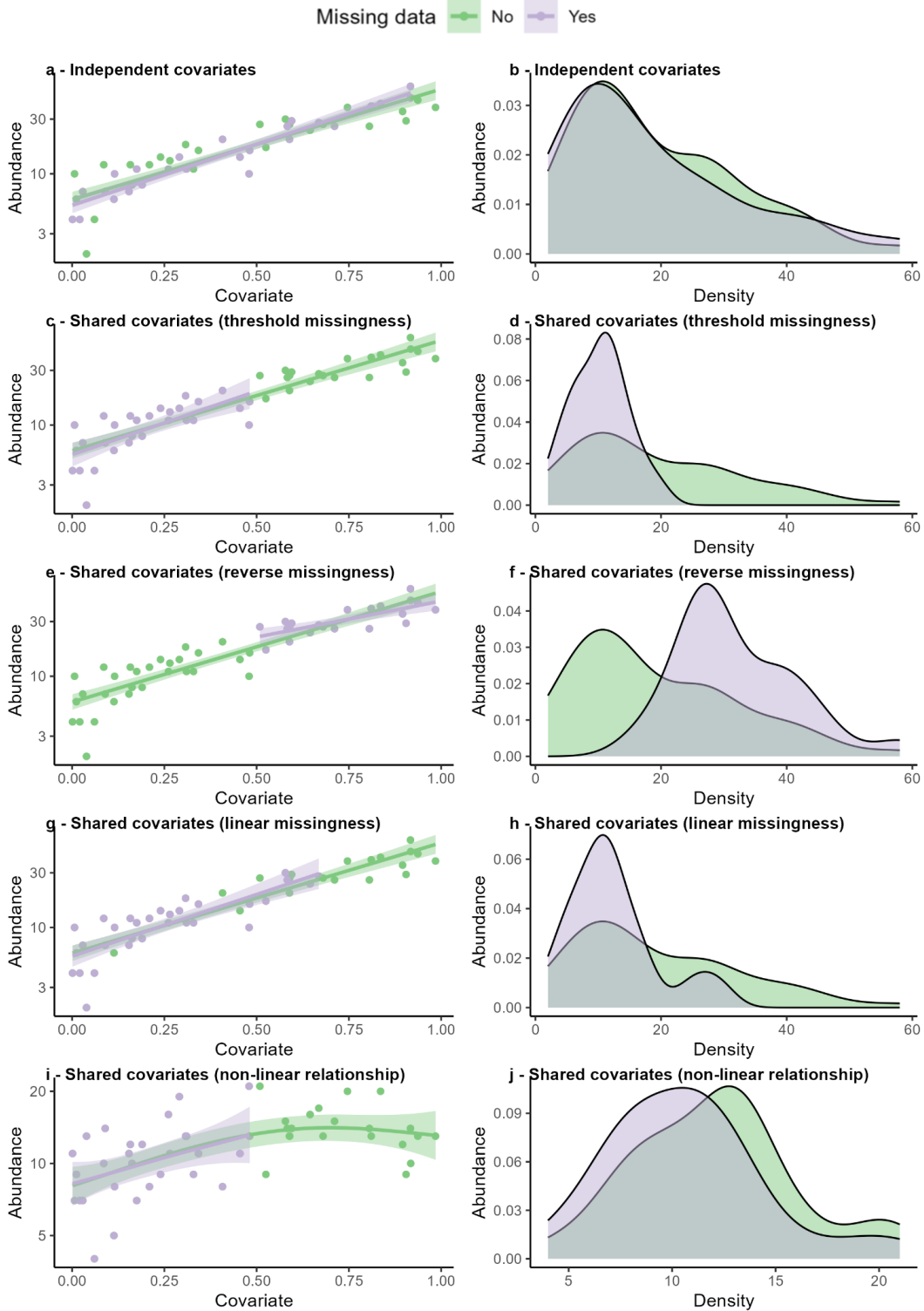252  data are presence/absence or presence-only (Baker *et al.*, 2022).

253       We begin considering the scenario when presence-absence data are available. In this case,

254  if there is a simple linear relationship, missing data do not necessarily cause bias in the estimated

255 effect of the covariate on biodiversity, even when missingness depends on the same covariate

256 (Fig. 1a, c, e, g) (Collins *et al.*, 2001). For instance, we could estimate the effect of elevation on

257 species occurrence, if it is linearly related, even if elevation is also associated with data gaps.

258 This is because the relationship between the covariate and species occurrence can be estimated

259 without bias using data over a restricted range of covariate values. This is shown in e.g., Fig. 1c -

260 the same relationship is found with a full dataset (green) or a restricted dataset with data gaps

261 (purple). Missing data can, however, cause problems when the underlying relationship between

262 the covariate and species occurrence is non-linear. In this case, data gaps can hinder estimating

263 the true form of the relationship (see Fig. 2i - a curved relationship in fit with the full dataset but

264 a simple positive linear relationship with the restricted dataset). The relationship that is fit using

265 the restricted dataset will critically depend on which portion of the covariate range is sampled.

266 Since many ecological associations show some non-linearity, or context-dependencies such that

267 relationships depend on the value of other variables (Spake *et al.*, 2023), we expect this issue is

268 likely to be widespread in species distribution models.

269       We now consider the alternative scenario of fitting a distribution model with presence-

270 only data. In this case, any data gaps could represent a lack of sampling or a lack of true species

271 occurrence. This creates an inherent identifiability challenge for any model seeking to separate

272 the processes affecting sampling from the true ecological processes affecting species

273 distributions with presence-only data (Hefley *et al.*, 2013; Baker *et al.*, 2022). Many methods

274 have been developed to generate pseudo-absences (Barbet-Massin *et al.*, 2012; Hertzog, Besnard

275 & Jay-Robert, 2014), but such models are still usually more prone to biases when there are

276 shared covariates affecting sampling and species occurrence (Baker *et al.*, 2022). More recent

14

277 approaches to modelling presence-only data, by integrating them with any available presence-

278 absence data (Fithian *et al.*, 2015), may help minimise some of these biases.

279



280

281  **Fig. 2 The impacts of different missing data mechanisms on regression (left) and sample**

282  **distributions (right).**

283  We use a hypothetical dataset to highlight different missing data mechanisms. In (a) and (b), the covariate

284  affecting sampling probability is *independent* from the covariate affecting species abundance. In this case,

285  both the estimated effect of the covariate (e.g., in a linear regression, shown in a by the solid line) and the

286  sample distribution (b) are similar in a dataset with (purple) and without (green) missing data. (i.e.,

287  missingness is MCAR). In (c) and (d), the covariate affecting sampling probability is the *same as or*

288  *correlated with* the covariate affecting species abundance - in this case, data are missing when the

289  covariate is above average (i.e., threshold missingness). The estimated effect of the covariate is the same

290  in the dataset with and without missing values (shown in c) but the sampling distribution is different (d).

291  In (e) and (f), the missingness pattern is reversed compared to (c) and (d)  (i.e., data are missing when the

292  covariate is below average), but we can similarly retrieve the same unbiased covariate effect (e) even

293  though there is greater mean abundance in the dataset with missing values (f). In (g) and (h), the covariate

294  affecting sampling probability is the *same as or correlated with* the covariate affecting species abundance

295  - in this case, the probability of missing data increases with the value of the covariate (i.e., linear

296  missingness). Again, the estimated effect of the covariate is the same (shown in g) but the sampling

297  distribution is different (h). In (i) and (j), the covariate affecting sampling probability is the *same as or*

298  *correlated with* the covariate affecting species abundance; additionally, the true relationship between the

299  covariate and species abundance is non-linear and data are missing when the covariate is above average.

300  The mechanism is now MNAR since the model cannot be correctly specified with the observed data.
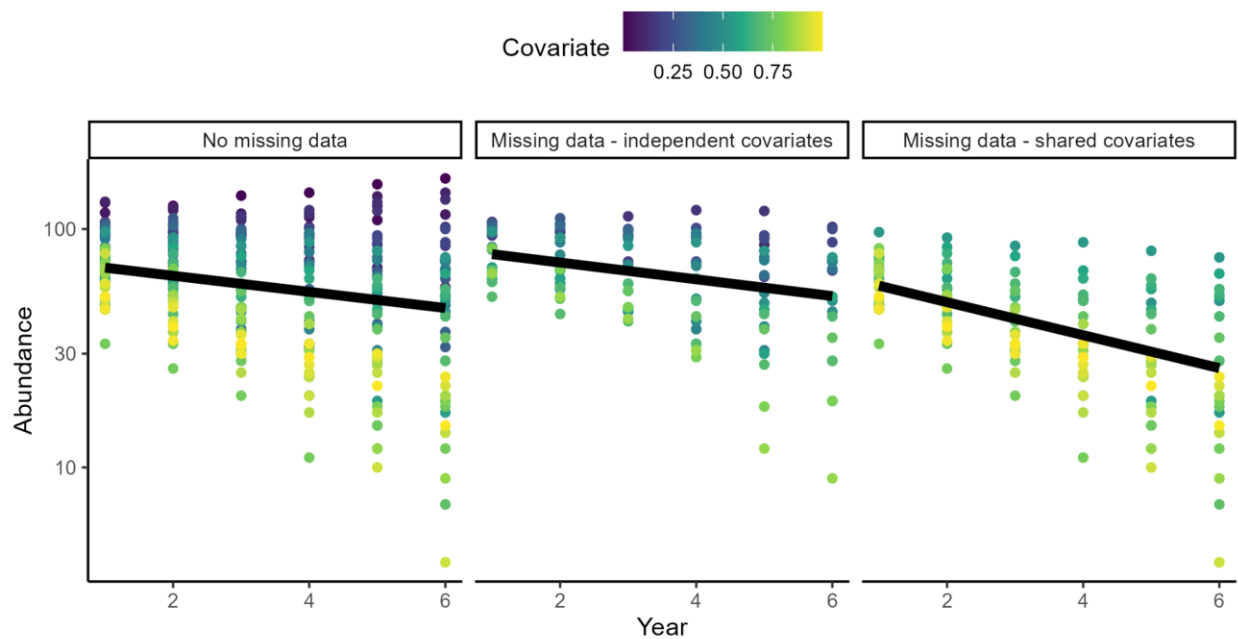
301

302  ## (2) Estimating trends in species abundances

303  Models to estimate species' trends tend to be descriptive: spatial variation is modelled by

304  including site identity (as a fixed or random term) while any temporal trend is modelled as a

305  simple year effect (either as a linear, spline or a categorical term) (Amano *et al.*, 2012; Bled *et*

16

306 *al.*, 2013). Drivers of the trend are not explicitly modelled when the goal is to simply estimate

307 the mean trend over time. Broader inferences about the trend estimated by such models are based

308 on the assumed representativeness of the sample, or prior knowledge of the inclusion

309 probabilities of sampling units (see design weights discussed in section IV 2). Basing inference

310 from the sampling design is the most traditional approach to surveys (Smith, 1976) and the

311 approach typically taken by official governmental surveys (van den Brakel & Bethlehem, 2008).

312 This approach has the advantage of avoiding complex assumptions in the statistical analysis

313 (Buckland *et al.*, 2012) and is perhaps also easier to analyse and communicate to stakeholders

314 and laypersons.

315     Simple trend models may, however, lead to biased trend estimates for biodiversity when

316 data gaps are not MCAR. We illustrate this in a simple simulation in which site-level species

317 trends were assumed to depend on a site-level covariate e.g., urban cover (Fig. 3). We assumed

318 sites were sampled either with a probability affected by an independent covariate (Fig. 3 middle

319 panel) or with a probability affected by the same site-level covariate affecting species trends

320 (Fig. 3 right panel), a scenario already identified as a pitfall in some monitoring schemes

321 (Buckland & Johnston, 2017). We estimated trends using a simple mixed effect model including

322 site and year. This shows that when the site-level covariate affected both sampling effort and

323 species' trends, the trends were biased, but site-level trends were unbiased when an independent

324 covariate affected sampling. In real world situations, many factors will influence the trend of a

325 species, but this toy simulation highlights the potential for bias caused by shared covariates.

326 Since the specific covariates affecting sampling effort and biodiversity trends are not considered

327 in the typical forms of analysis for trend modelling, trend analyses are liable to be affected by

328 MNAR, whereas by including appropriate covariates (where possible), the data gaps become

329  MAR instead and trends will be unbiased. Without conditioning on the covariates involved, trend

330  estimates might be overestimated if missing data are more common in static regions where

331  species trends are more stable; but underestimated if missing data are more common in dynamic

332  regions where species trends more strongly deviate from zero (Fig. 3) (Bowler *et al.*, 2022;

333  Buckland & Johnston, 2017).



334

**Fig. 3 The impacts of different missing data mechanisms on trend modelling**

336  We use a hypothetical scenario in which a mean trend model is fit to datasets that vary in their missing

337  data mechanism. We assumed a scenario of 50 sites that varied in an environmental covariate affecting

338  species trends (trends were stable or even increasing at low values of the covariate and declining at

339  increasingly high values of the covariate). When missing data was independent (i.e., a MCAR pattern -

340  the covariate affecting sampling probability was a different and uncorrelated covariate), the overall mean

341  trend (estimated by the year effect in a generalized linear mixed effect model that also included a site

342  random effect) was similar with (middle panel) and without (left panel) missing data. By contrast, when

343  the same covariate affected both species' trends and sampling probability, leading to less sampling in sites

344  with low values of the covariate (notice there are fewer blue points in the right panel - a MNAR pattern),

345  the overall mean trend was downward biased with missing data (right panel) compared to the scenario of

346  no missing data (left panel).


# IV. Missing data solutions

348  A broad range of methods to deal with missing data have been used in ecology (Hossie, Gobin &

349  Murray, 2021; Nakagawa & Freckleton, 2008; Lopucki *et al.*, 2022). Many solutions are

350  particularly relevant when data are missing in both response and predictor variables. Here, we

351  focus on the typical scenario in biodiversity modelling of missing data only in the response

352  variable (i.e., in the biodiversity data) since typical predictors tend to have no or few gaps (e.g.,

353  site identity or environmental data from remote sensing). We organise solutions into three groups

354  - subsampling, weighting and imputation (Fig. 4), which have been tested to varying degrees

355  already with both structured and unstructured biodiversity data (Table 2). Most solutions to deal

356  with missing data are only appropriate for MCAR or MAR missingness. MNAR is the most

357  challenging missing data mechanism to deal with in statistical modelling so, we deal with

358  MNAR in a later section.

359

**Fig. 4 Visualisation of contrasting approaches to deal with data gaps.**

We focus on spatial gaps to illustrate the possible approaches, but the ideas apply to other types of data gaps (Fig. 1). (top) the landscape is divided into four quarters (e.g., representing different habitats or geographic regions). One quarter (top right quarter) has been sampled more (4 sampling sites) than the others (2 sampling sites). Solutions: Random subsampling (bottom left): two sites are randomly subsampled from the oversampled quarter to create a dataset with an even sampling coverage across quarters. Weighting (bottom middle): data from the oversampled quarter is downweighted in the statistical model so data from all quarters similarly influence the modelled results. Imputation (bottom right): missing values at unsampled sites are imputed based on the spatial pattern in the data and/or environmental covariates, and summary parameters are calculated based on both predictions at sampled and unsampled sites. In subsampling and weights, the aim is to improve the representativeness of the sample for statistical inference at the population-level. In imputation, the aim is to directly predict population-level values.

373

374 **Table 2 Example applications of the solutions to deal with data gaps within biodiversity data.**

| Type of data gaps | Typical approaches: |
|---|---|
| Within-year | Sometimes imputed e.g., spline terms to smooth over seasonal variation in sampling times during the flight period of butterflies (Dennis *et al.*, 2016) |
| Annual | Sometimes imputed e.g., general linear models to impute annual gaps based on mean site and year effects, optionally allowing for habitat effects e.g., used in TRIM abundance indices, (Lehikoinen *et al.*, 2016) |
| Spatial | Often ignored, but occasionally weighting by geographic regions (Bled *et al.*, 2013) or imputed (Breivik *et al.*, 2021),  or reduced by subsampling (Johnston *et al.*, 2021). |

375

## 376 **(1) Subsampling**

377 The 'Big Data Paradox' highlights that there can be trade-offs between dataset size and dataset

378 quality (Bradley *et al.*, 2021; Meng, 2018). Small datasets can be preferable to large datasets, if

379 they are more representative and less heterogeneous than a larger dataset (Bayraktarov *et al.*,

380 2019). Based on such thinking, some studies have proposed to 'reverse engineer' structure in

381 biodiversity data by filtering data (Rapacciuolo, Young & Johnson, 2021). Part of this reverse

382 engineering has attempted to deal with spatial biases; for instance, by spatially subsampling data

383 to reduce the unevenness of sampling effort across the landscape (Steen *et al.*, 2021; Matutini *et*

384 *al.*, 2021; Steen, Elphick & Tingley, 2019; Boria *et al.*, 2014; Robinson *et al.*, 2020). This has

385 been tested on, for instance, the semi-structured data compiled by eBird (Johnston *et al.*, 2020).

386 Some have also applied this approach to reduce temporal skews in sampling effort (Hof &

387 Bright, 2016; Zbinden *et al.*, 2014), although not always successfully (Callcutt, Croft & Smith,

388 2018). Subsampling can also be used to balance the amount of data across a single or multi-

389 dimensional environmental gradient; essentially stratified sampling of the original sample

390  (Meng, 2022; Nunez-Penichet *et al.*, 2022). Recent class balancing approaches have been

391  developed to ensure that important observations, especially for rare species, are not lost during

392  the subsampling process (Robinson *et al.*, 2020; Steen *et al.*, 2021; Gaul *et al.*, 2022).

393

## (2) Weighting

395  Weighting is a common practice in survey analysis, especially in the social sciences (Li *et al.*,

396  2013; Seaman & White, 2013; Raghunathan, 2004). Weighting can serve different purposes,

397  including reducing the impact of confounding variables when the goal is to estimate the causal

398  effect of an intervention. But weighting can also be used to deal with missing data that is not

399  MCAR. For instance, weighting can be used to reduce selection bias caused by participant

400  nonresponse in surveys (Seaman & White, 2013), but it is less often used to account for data

401  gaps in biodiversity data (Boyd, Powney & Pescott, 2023a; Aubry & Francesiaz, 2022).

402  Different types of weights have been used in the analysis of biodiversity data: (1) design

403  weights; (2) estimated non-response weights (or sampling weights) and (3) population weights.

404  Each form of weighing is intended to improve sample representativeness of some target

405  population but vary in terms of whether the weights derive from the sampling design and the

406  dimension of representativeness under consideration. Design weights are based on the study

407  sampling design and assumed to be known with certainty, and hence are only relevant for

408  structured schemes with a sampling design. For instance, in many national bird breeding

409  schemes, the design weights are based on the geographic strata that underlie a random stratified

410  study design (Buckland *et al.*, 2012). Non-response weights are used to account for unplanned

411  missing data in structured schemes (Frair *et al.*, 2004) or variation in sampling effort in

412  unstructured schemes (Johnston *et al.*, 2020; Hefley *et al.*, 2013), which means that are not

413    known with certainty and must be estimated. Population weights are used to ensure the sample is

414    representative of the full distribution/population of a species and are typically assumed to be

415    known. Population weights are used in the calculation of supranational/international indicators

416    (e.g., farmland or woodland bird indicators (Gregory *et al.*, 2005)) in which national estimates

417    are combined by giving greater weight to regions that harbour a larger proportion of the species

418    total population.

419         Non-response weights are usually the most difficult to include since they are not known *a*

420    *priori* and need to be estimated. Predictive models (e.g., random forest models) have been used

421    to predict the probability that a site is sampled based on the set of covariates (e.g., land cover or

422    climate, or accessibility) available across all sampled and unsampled sites, with the inverse of

423    these probabilities used as weights (Little *et al.*, 2022; Johnston *et al.*, 2020). Alternatively,

424    poststratification (for categorical covariates), or more generalized calibration approaches

425    (allowing both continuous and categorical covariates), can be used, which adjust the weight

426    given to each data point until the joint or marginal distributions of covariate values in the

427    observed sample matches those for the population (Boyd, Stewart & Pescott, 2023b). In both

428    cases, weighting can cause problems when there are regions within the target population with

429    close to zero probability of being sampled, which could lead to some data points having

430    extremely large weights. In this case, weights may need to be redefined e.g., by coarsening the

431    covariates used to define the weights, or by truncating weight values so that extreme weights are

432    not produced (Battaglia, Hoaglin & Frankel, 2009). Poststratification can be preceded by multi-

433    level regression (for so-called "Mr P" analysis) for partial pooling of information across strata

434    before poststratification of the model predictions, which may be especially useful when some

435    strata contain few data points (Gelman, 2007).

436    The most appropriate approach is likely to be question- and taxon-specific, varying with

437    how much the species range extends across the region of interest. For example, it would usually

438    not be important to upweight under-sampled regions where a species is rare, or even absent,

439    when estimating trends in its total population size. If, however, the goal is to estimate trends in

440    the average site-level population trend, then it would be important to up-weight data from under-

441    sampled regions, even from where the species is rare. For instance, in the UK bat monitoring

442    scheme, data are weighted to allow for the different sampling rates across England, Scotland and

443    Wales in proportion to the ratio of non-upland area to number of sites surveyed for the relevant

444    country (Bat Conservation Trust, 2023). However, this weighting is not applied to range

445    restricted species, such as the serotine bat, *Eptesicus serotinus* that is only found in southern

446    England.

447

448    **(3) Imputation**

449    Imputation involves replacing missing values in a dataset with plausible estimates. A range of

450    imputation procedures have been developed, which can fill gaps in both response and predictor

451    variables (Carpenter & Kenward, 2012). Imputation is probably the most flexible and widely

452    used approach to account for missing data across ecology and beyond. In biodiversity modelling,

453    missing values are more often concentrated in the response variable (i.e., the biodiversity value),

454    hence imputation here can be equated with making model predictions at unsampled sites and

455    times.

456    Imputation is already in use in biodiversity trend monitoring, especially to account for

457    within-year and annual data gaps (Table 3). Early approaches used chain indices or route

458    regression (Ter Braak *et al.*, 1992) or the Underhill index, using an expectation-maximisation

459 algorithm (Underhill & Prysjones, 1994) designed for waterbirds (Rehfisch *et al.*, 2003). A range

460 of further model-based approaches have been developed that fill data gaps using mean effects of

461 site and year, e.g., to fill annual gaps using TRIM/birdSTATs, commonly used for bird indices

462 (Lehikoinen *et al.*, 2016); or using splines e.g., to fill seasonal gaps in butterfly data (Schmucki

463 *et al.*, 2016; Dennis *et al.*, 2016) or using ecological covariates (Dakki *et al.*, 2021). A Bayesian

464 framework is especially useful for dealing with missing values in the response since they are

465 naturally imputed with a full probability distribution during model fitting. Bayesian occupancy-

466 detection models have been used to analyse opportunistic species observations from citizen

467 science, with annual data gaps imputed before the predicted annual proportion of occupied sites

468 is calculated (Outhwaite *et al.*, 2019). The flexibility of Bayesian models means they could also

469 incorporate expert knowledge as priors as a way to help fill data gaps (Johnson *et al.*, 2023).

470  While imputation is already used to deal with annual and within-year gaps, it has been

471 less often used to deal with spatial gaps when the focus is mean trend modelling of species'

472 abundances or occurrences. An exception is studies of changes in species' range sizes, which use

473 distribution models to predict the full distribution of a species at multiple time points, before

474 change is assessed (Grattarola, Bowler & Keil, 2023). Monitoring schemes with large spatial

475 coverage are also beginning to use distribution or abundance models to predict spatio-temporal

476 patterns of abundance change across whole countries (e.g., eBird maps and BTO maps). In these

477 cases, regression models fit to the available data make predictions at unsampled sites based on

478 the effects of environmental covariates and/or spatial structure (Bush *et al.*, 2017; Ver Hoef *et*

479 *al.*, 2021; Breivik *et al.*, 2021). Geostatistical methods also offer a range of interpolation

480 methods for spatial data, including kriging, which are especially useful when there is a strong

481 spatial pattern in the data (Ballesteros-Mejia *et al.*, 2013; Kreft & Jetz, 2007; Lin *et al.*, 2008).

# V. Pro and cons of each solution

**Table 3 Summary of the pros and cons of each approach to deal with missing data in biodiversity monitoring**

| Solution | Pros | Cons |
|---|---|---|
| Subsampling | - arguably the simplest approach, especially for spatial gaps<br><br>- already a routine feature of many species distribution modelling protocols<br><br>- aligns with rarefaction approaches used in community ecology | - could mean excluding a large amount of data, which may be unacceptable for citizen science and engaging/retaining volunteers<br><br>- most protocols focus on a single dimension (e.g., filtering by geographic region)<br><br>- more complex to implement when gaps are multi-dimensional or temporally varying |
| Weighting | - standard practise to deal with sample unrepresentativeness in other disciplines, especially social sciences | - poorly understood in ecology<br><br>- diverse range of possible weighting techniques (Valliant, 2020; Boyd *et al.*, 2023b) but little ecological guidance available to help selection |
| Imputation | - suitable approach if missing data are within the environmental covariates as well as within the biodiversity response<br>- offers the promise to generate the continuous space-time data cubes of the Essential Biodiversity Variable framework (Kissling *et al.*, 2018; Jetz *et al.*, 2019). | - becomes inefficient as missingness increases, e.g., when the number of unsampled locations/times is large<br><br>- requires a good understanding of the ecological system to predict the missing biodiversity values |

All of the approaches have the potential to reduce the bias in parameter estimates but differ in complexity, scope and typical practice (Table 3) (Little *et al.*, 2022; Collins *et al.*,

488   2001). Moreover, while we separated the methods into three categories for convenience, their

489   distinctions are not absolute. For instance, subsampling essentially assigns those population units

490   included in the subsample a weight of 1 and the remainder a weight of 0. Often, but not always,

491   the reduction in bias due to application of the above solutions comes at a cost of increasing

492   parameter uncertainty: the classic bias-variance trade-off (Hefley *et al.*, 2013). This is because

493   subsampling directly reduces the sample size; weighting reduces the effective sample size; and

494   imputation adds uncertainties via predictions at unsampled points. But this trade-off does not

495   always apply; for instance, poststratification can lead to the dual benefits of reduced bias and

496   increased precision depending on the choice of covariates (Little & Vartivarian, 2005).

497         Covariates used to account for data gaps are often called 'auxiliary variables' (Little *et*

498   *al.*, 2022), which are typically not of central interest to the scientific questions but are included in

499   one or more of the analysis steps for subsampling, weighting or imputing. The general

500   recommendation from the missing data theory and survey sampling literature is to be generous

501   when deciding which covariates to use to adjust for data gaps, considering covariates relating to

502   the missingness (i.e., sampling effort in the context of biodiversity data gaps) to reduce bias and

503   those related to the biodiversity outcome to reduce the variance (Collins *et al.*, 2001; Caughey *et*

504   *al.*, 2020). It is worth noting, however, that selecting auxiliary variables on a purely correlative

505   basis can increase bias in some circumstances (Thoemmes & Rose, 2014), and a safer strategy is

506   to select them on theoretical grounds (Mohan & Pearl, 2021). When auxiliary variables are

507   related to both the biodiversity outcome and the pattern of missingness, weighting approaches

508   can reduce bias and improve precision (Little & Vartivarian, 2005). The success of any of the

509   solutions, hence, critically depends on the choice of auxiliary variables (Little *et al.*, 2022). A

510   recent study testing the use of weighting approaches to account for spatial biases in a reasonably
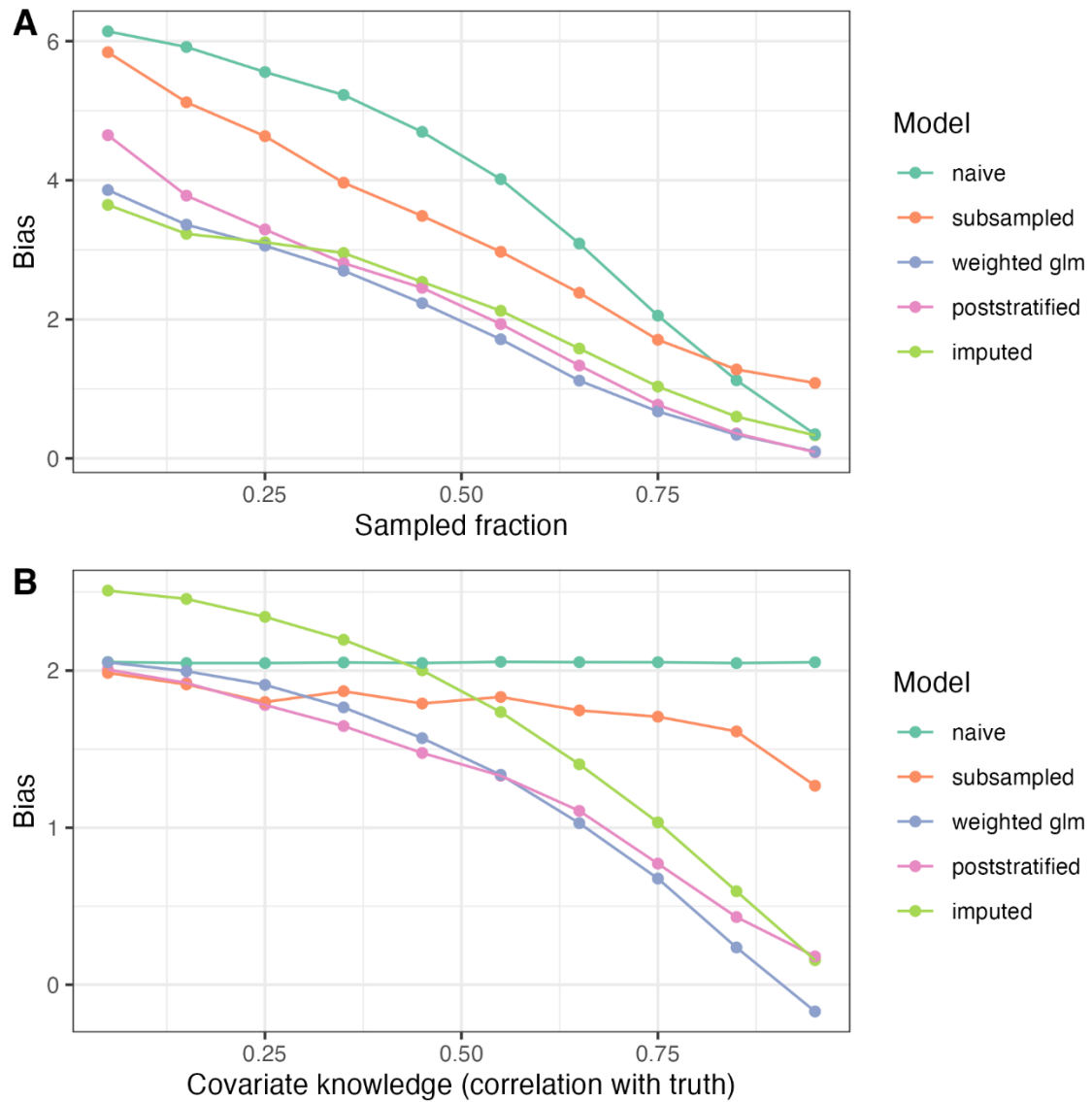
511    well-understood ecological system found that the selected auxiliary variables had only limited

512    success in mitigating bias (Boyd *et al.*, 2023b).

513    We illustrate some of these challenges and the application of each potential solution with

514    a toy example of an abundance dataset with missing values (Fig. 5). We simulated a landscape in

515    which a covariate (let's say representing 'habitat quality') affected both species abundance and

516    the likelihood of a site being sampled. The analysis aimed to estimate the mean abundance of the

517    species across all sites in the landscape. We varied the total fraction of sites that were sampled

518    and the degree of knowledge available on the covariate affecting sampling/species (modelled as

519    the correlation between the covariate involved in the data generation process and the covariate

520    available to the modeller). We compared subsampling, weighting and imputation, which all used

521    the available covariate data for adjustment. For subsampling, we subsampled one site at random

522    at each habitat quality value. For weighing, we compared two approaches: fitting a weighted

523    regression model using model-robust sandwich variance estimators or using a poststratification

524    approach. For imputation, we fit a Bayesian model using JAGS in which NA values were

525    inserted to represent the missing response data.

526    The results show that all methods do better than a naive approach that did not attempt to

527    account for missingness in the estimation of the mean abundance (Fig. 5). Subsampling

528    performed the worst, while weighting and imputation performed similarly. Poststratification

529    tends to perform less well with a lower sampling fraction i.e., when the number of missing

530    values was high (Fig. 5A), because the sample did not always contain all the habitat quality

531    values found in the population and the weighing could not account for entirely unsampled

532    regions. All models performed less well at the available covariate became a weaker proxy of the

533    true driving covariate (Fig. 5B). In further simulations, we found that imputation performed

534 poorer when there were additional covariates affecting species abundance and these covariates

535 were not modelled, highlighting the importance of understanding the ecological system for

536 imputation (Fig. S1). We do not intend this simulation to be exhaustive - rather to highlight the

537 potential ways in which the availability of data and degree of knowledge about the factors

538 causing bias and the availability of covariate data affects any attempts to account for missing

539 data.

540 We point the reader towards some useful R packages and functions in the Supporting

541 Information (Table S1).

**542**

**543**     **Fig. 5 The ability of missing data solutions to adjust for bias in biodiversity data.**

**544**      We assumed a landscape of 400 cells and that a covariate affected both species abundance and the likelihood of a

**545**     cell being sampled. In A: we vary the fraction of the cells that were sampled. In B: we vary the correlation between

**546**     the true covariate and the covariate available for analysis, as measure of the available knowledge (correlation of 1 =

**547**     perfect covariate and knowledge). The models to estimate the parameter of interest (mean abundance) were: naive

**548**     (no correction); subsampled (cells were subsampled along the covariate gradient), weighted (two methods: weighted

**549**     glm using the svyglm function, and weighted by poststratification, using postStratify, both in the survey package)

550    and imputed (using JAGS to impute NAs in the response). Points show the mean bias (difference between model

551    prediction and truth) across 100 independent runs.

552

# VI. Dealing with Missing Not at Random

554    Dealing with Missing Not at Random (MNAR) is more challenging than dealing with the other

555    data mechanisms (Little & Rubin, 2019). In this case, missingness is directly associated with

556    unavailable data, which could be either the missing biodiversity values or missing covariate data

557    that are not known to be important or are not measured/measurable. This makes MNAR

558    especially difficult to diagnose (but see Conn et al. (2017) for suggestions) and model, since

559    auxiliary variables are not available. MNAR can arise through a number of mechanisms in

560    biodiversity monitoring data.

561        MNAR can be an outcome of preferential sampling - more intense sampling effort where

562    the species is expected (Diggle, Menezes & Su, 2010; McClure & Rolek, 2023) - which leads to

563    more missing values in places where the species is rare or absent. Preferential sampling can

564    arise, for instance, if observers visit a location to specifically observe a species that others have

565    observed there before (Laney *et al.*, 2021; Pennino *et al.*, 2019). Preferential sampling can also

566    be a planned sampling strategy (Alessi *et al.*, 2023). For rare species, preferential sampling can

567    be optimal when the goal is to estimate species detection probability and account for imperfect

568    detection, since sufficient observations of the species can only be achieved by sampling where

569    they are more common (Specht *et al.*, 2017). Similarly, it can be optimal to expend greater

570    sampling effort where the species is common if the goal is to estimate trends in the total

571    population size, since regions where the species is scarce are less important for the overall trend.

572    For organisms associated with specific habitats, such as wetland species or colonial seabirds,

573    dedicated structured monitoring schemes target their habitats (McClure & Rolek, 2023). In such

574    schemes, missing data outside of these core habitats are not considered part of the target

575    population.

576        Typical approaches to modelling data allowing for MNAR are selection models

577    (Heckman, 1979) and pattern-mixture models (Herzog and Rubin, 1983). Both model the joint

578    distribution of the data and the data availability, but differ in how these processes are

579    decomposed. Both also require making strong assumptions about the missing data mechanism,

580    but can be useful to explore the consequences of plausible options as a sensitivity analysis

581    (Little, 1995). In the ecological literature, preferential sampling has been modelled using marked

582    point process models, which jointly model the sampling intensity (the points), the biodiversity

583    value at those points (the marks) and the dependence between them (Conn *et al.*, 2017; Pennino

584    *et al.*, 2019; Laxton *et al.*, 2023). Another approach to inference in a NMAR scenario is to use

585    instrumental variables i.e., variables that affect the probability of sampling/data availability but

586    are independent of the biodiversity variable of interest (Tchetgen & Wirth, 2017; Bailey, 2023).

587    The challenge, however, is to identify such variables.

# VII. General guidelines for dealing with biodiversity data gaps

590    Our review highlights the potential value of 'missing data thinking' when analysing biodiversity

591    data. We argue that MCAR data gaps are unlikely in most biodiversity data contexts, which

592    means that researchers will need to consider whether and how they deal with data gaps in their

593    analysis. While it is premature to make very specific guidelines, we summarise here some of the

594    considerations needed when dealing with data gaps in biodiversity data at different stages of data

595    collection, analysis and reporting.


596    **(1) Study design**

597    For new monitoring schemes, planned data gaps that deviate from MCAR (i.e., a random

598    sample) can be seen as opportunities rather than challenges since solutions are available to deal

599    with missing data. Intentionally missing some data has been proposed for ethical or practical

600    reasons in some study designs e.g., (Noble & Nakagawa, 2021; Herrera, 2019). In citizen

601    science, planned data gaps could help increase uptake and avoid participant fatigue, especially

602    caused by collecting difficult data. For instance, the UK Breeding Bird Survey includes an

603    'upland rovers' component in which the standard protocol is modified to allow for fewer visits to

604    remote sites (Darvill *et al.*, 2020). Alternative study designs, such as wave missingness (Little &

605    Rhemtulla, 2013) or a rotating panel design (Nielsen *et al.*, 2009) may increase the sustainability

606    of long-term monitoring for some taxa or regions with few willing participants. But such an

607    approach has to balance the cost of increased study design complexity and potential implications

608    for the range of questions that can be addressed.


609        For existing monitoring schemes, data gaps may be filled, where possible, by promoting

610    data collection in certain areas. Within citizen science projects, there is evidence that participants

611    can be nudged to collect more data in regions identified as sampling priorities (Callaghan *et al.*,

612    2023; Callaghan *et al.*, 2019). Previous studies have identified sampling priorities in different

613    ways; for instance, based on the expected influence of a data point (Callaghan *et al.*, 2019) or

614    predictions based on species distribution models (Chiffard *et al.*, 2020). Since data collected by

615    monitoring schemes are often collected for multiple purposes, the challenge is identifying the

616    common set of sampling priorities.

617    For synthesis studies compiling data from independent studies, data mobilisation efforts

618    may be tailored to improve sample representativeness of the target population, by expending

619    more effort to under-sampled units. This could be informed by exploring the transferability of

620    model predictions across spatial or temporal units based on currently available data (Spake *et al.*,

621    2022). Regions with high transferability may represent appropriate sampling strata to guide

622    mobilisation efforts. Moreover, these sampling strata may inform the adjustment for data gaps in

623    subsequent modelling of the population mean.

624    **(2) Evaluating and reporting missingness**

625    Developing a causal model (e.g., using a DAG) of the factors affecting sampling probability and

626    biodiversity can be useful first step to identify auxiliary variables for adjusting data gaps –

627    variables linked to both sampling probability and biodiversity are those creating bias (Mohan &

628    Pearl, 2021). As far as possible, data should then be collected on the covariates that are likely to

629    explain missingness. Statistical models can be used to test whether covariates that are associated

630    with missingness are also associated with biodiversity patterns, though of course this is only

631    possible in the sampled data. Unplanned missingness in structured schemes could be investigated

632    by disseminating follow-up surveys to participants to determine their reasons for missed surveys.

633    Follow-on data collection, e.g., with paid surveys, in regions or times of missing data could also

634    help understand whether there are fundamental differences in biodiversity patterns between the

635    original dataset and the extended dataset.

636    Missingness, and how it is dealt with, tends to be insufficiently reported in biodiversity

637    trend analyses. Some reporting frameworks for missing data have been developed for other

638    disciplines (Lee *et al.*, 2021) but are in their early stages in ecology (Boyd *et al.*, 2022). At a

639    minimum, we propose that missingness can be reported in terms of the proportion of sampling

640    units that are spatial, annual and within-year gaps, and the number of unplanned gaps for

641    structured monitoring schemes (Fig. 1). Visualizations of the distributions of covariates in

642    sampled and non-sampled times/sites could also effectively highlight key systematic differences.


643    **(3) Modelling to account for data gaps**

644    The impact of data gaps depends on multiple factors: whether the factors affecting missingness

645    are independent of the factors affecting biodiversity and biodiversity itself; the ecological

646    questions being asked and which covariates are available and included in the analysis. Because

647    of this, potential impacts of missingness have to be considered for each species-question-dataset

648    combination. A dataset *per se* is not biased. Subsampling, weighting and imputation all have the

649    potential to reduce bias caused by data gaps. Many, but not all, solutions will navigate the bias-

650    variance trade-off. Weighting is probably the most under-used in ecology and could be applied

651    more often, especially to account for spatial gaps when the goal is estimating mean abundance or

652    abundance trends. Imputation methods offer the potential to fill in spatio-temporal gaps to

653    generate the space-time data cubes of the Essential Biodiversity Framework (Kissling *et al.*,

654    2018), but its success is dependent on the ability to model the variation in the biodiversity

655    response. Since available covariates are likely to be only partly successful in reducing bias,

656    sensitivity analysis could be help explore how different assumptions of missingness change the

657    model coefficients and predictions, and the uncertainties of them (Little, 1995; Leurent *et al.*,

658    2018). For some contexts, it might be more statistically efficient and ecologically interpretable to

659    redefine the target region of interest to a region with fewer data gaps.

# VIII. Conclusions

(1) Biodiversity datasets containing information on species' occurrences and abundances are rapidly growing in size, but data gaps are not necessarily closing. Nonetheless, big biodiversity datasets are invaluable for a broad range of basic and applied questions, and increasingly for policy-relevant questions about the status and trends of biodiversity at large-scales. Heterogeneity in sampling efforts - whether by citizen scientists or scientists - creates different types of data gaps in the available data. Such data gaps are among the biggest hindrances to making use of these growing data sources for large-scale inferences.

(2) We show how 'missing data thinking' can help decide whether a data gap is problematic in a given context and provides directions on possible solutions. We show that an important determinant of bias is whether factors affecting sampling effort are correlated with those affecting biodiversity: shared covariates affecting sampling effort and biodiversity have the potential to lead to biased analyses if not taken into account.

(3) Multiple approaches are available to account for missing data but they depend on knowledge and availability of covariates associated with missingness. A lack of training for ecologists in commonly employed approaches in other disciplines has meant there are few standard practices in ecology to deal with gaps. We highlight multiple methods that are ripe for comparison across different ecological problems.

(4) At the same, statistical solutions can only go so far, closing data gaps with more coordinated data collection across monitoring stakeholders is also important as we move forwards.

# IX. Acknowledgements

# X. Supporting Information

**Table S1** Selected R tools that can help with missing data problems and their potential

application for use in biodiversity research.

**Fig. S1** Extended analysis of missing data solutions when additional covariates affect the

biodiversity response.

R script for the example solution simulations (Fig. 5) are here:

https://github.com/bowlerbear/dataGaps

37

# IX. References

ALBERT, C. H., YOCCOZ, N. G., EDWARDS, T. C., GRAHAM, C. H., ZIMMERMANN, N. E. & THUILLER, W. (2010). Sampling in ecology and evolution - bridging the gap between theory and practice. *Ecography* **33**(6), 1028-1037.

ALESSI, N., BONARI, G., ZANNINI, P., JIMENEZ-ALFARO, B., AGRILLO, E., ATTORRE, F., CANULLO, R., CASELLA, L., CERVELLINI, M., CHELLI, S., DI MUSCIANO, M., GUARINO, R., MARTELLOS, S., MASSIMI, M., VENANZONI, R., ZERBE, S. & CHIARUCCI, A. (2023). Probabilistic and preferential sampling approaches offer integrated perspectives of Italian forest diversity. *Journal of Vegetation Science* **34**(1).

AMANO, T., LAMMING, J. D. L. & SUTHERLAND, W. J. (2016). Spatial Gaps in Global Biodiversity Information and the Role of Citizen Science. *Bioscience* **66**(5), 393-400.

AMANO, T., OKAMURA, H., CARRIZO, S. F. & SUTHERLAND, W. J. (2012). Hierarchical models for smoothed population indices: The importance of considering variations in trends of count data among sites. *Ecological Indicators* **13**(1), 243-252.

AUBRY, P. & FRANCESIAZ, C. (2022). On comparing design-based estimation versus model-based prediction to assess the abundance of biological populations. *Ecological Indicators* **144**.

BAILEY, M. A. (2023). A New Paradigm for Polling. *Harvard Data Science Review* **5**(3).

BAKER, D. J., MACLEAN, I. M. D., GOODALL, M. & GASTON, K. J. (2022). Correlations between spatial sampling biases and environmental niches affect species distribution models. *Global Ecology and Biogeography* **31**(6), 1038-1050.

BALLESTEROS-MEJIA, L., KITCHING, I. J., JETZ, W., NAGEL, P. & BECK, J. (2013). Mapping the biodiversity of tropical insects: species richness and inventory completeness of African sphingid moths. *Global Ecology and Biogeography* **22**(5), 586-595.

BARBET-MASSIN, M., JIGUET, F., ALBERT, C. H. & THUILLER, W. (2012). Selecting pseudo-absences for species distribution models: how, where and how many? *Methods in Ecology and Evolution* **3**(2), 327-338.

BARTLETT, J. W., HAREL, O. & CARPENTER, J. R. (2015). Asymptotically Unbiased Estimation of Exposure Odds Ratios in Complete Records Logistic Regression. *American Journal of Epidemiology* **182**(8), 730-736.

BAT CONSERVATION TRUST. (2023). The National Bat Monitoring Programme Annual Report 2022. Bat Conservation Trust, London. Available at [www.bats.org.uk/our-work/national-bat-monitoringprogramme/reports/nbmp-annual-report](www.bats.org.uk/our-work/national-bat-monitoringprogramme/reports/nbmp-annual-report).

BATTAGLIA, M. P., HOAGLIN, D. C. & FRANKEL, M. R. (2009). Practical Considerations in Raking Survey Data. *Survey Practice* **2**(5).

BAYRAKTAROV, E., EHMKE, G., O'CONNOR, J., BURNS, E. L., NGUYEN, H. A., MCRAE, L., POSSINGHAM, H. P. & LINDENMAYER, D. B. (2019). Do Big Unstructured Biodiversity Data Mean More Knowledge? *Frontiers in Ecology and Evolution* **6**.

BINLEY, A. D. & BENNETT, J. R. (2023). The data double standard. *Methods in Ecology and Evolution* **14**(6), 1389-1397.

BLED, F., SAUER, J., PARDIECK, K., DOHERTY, P. & ROYLE, J. A. (2013). Modeling Trends from North American Breeding Bird Survey Data: A Spatially Explicit Approach. *Plos One* **8**(12).

BOAKES, E. H., MCGOWAN, P. J. K., FULLER, R. A., DING, C. Q., CLARK, N. E., O'CONNOR, K. & MACE, G. M. (2010). Distorted Views of Biodiversity: Spatial and Temporal Bias in Species Occurrence Data. *Plos Biology* **8**(6).

BORIA, R. A., OLSON, L. E., GOODMAN, S. M. & ANDERSON, R. P. (2014). Spatial filtering to reduce sampling bias can improve the performance of ecological niche models. *Ecological Modelling* **275**, 73-77.

744 BOWLER, D. E., CALLAGHAN, C. T., BHANDARI, N., HENLE, K., BARTH, M. B., KOPPITZ, C., KLENKE, R., WINTER, M.,
745      JANSEN, F., BRUELHEIDE, H. & BONN, A. (2022). Temporal trends in the spatial bias of species
746      occurrence records. *Ecography* **2022**(8).
747 BOYD, R. J., POWNEY, G. D., BURNS, F., DANET, A., DUCHENNE, F., GRAINGER, M. J., JARVIS, S. G., MARTIN, G.,
748      NILSEN, E. B., PORCHER, E., STEWART, G. B., WILSON, O. J. & PESCOTT, O. L. (2022). ROBITT: A tool for
749      assessing the risk-of-bias in studies of temporal trends in ecology. *Methods in Ecology and*
750      *Evolution* **13**(7), 1497-1507.
751 BOYD, R. J., POWNEY, G. D. & PESCOTT, O. L. (2023a). We need to talk about nonprobability samples. *Trends*
752      *in Ecology & Evolution* **38**(6), 521-531.
753 BOYD, R. J., STEWART, G. B. & PESCOTT, O. L. (2023b). Descriptive inference using large, unrepresentative
754      nonprobability samples: An 1introduction for ecologists. *EcoEvoRxiv*.
755 BRADLEY, V. C., KURIWAKI, S., ISAKOV, M., SEJDINOVIC, D., MENG, X. L. & FLAXMAN, S. (2021). Unrepresentative
756      big surveys significantly overestimated US vaccine uptake. *Nature* **600**(7890), 695-+.
757 BREIVIK, O. N., AANES, F., SOVIK, G., AGLEN, A., MEHL, S. & JOHNSEN, E. (2021). Predicting abundance indices in
758      areas without coverage with a latent spatio-temporal Gaussian model. *Ices Journal of Marine*
759      *Science* **78**(6), 2031-2042.
760 BUCKLAND, S. T., BAILLIE, S. R., DICK, J. M., ELSTON, D. A., MAGURRAN, A. E., SCOTT, E. M., SMITH, R. I., SOMERFIELD,
761      P. J., STUDENY, A. C. & WATT, A. (2012). How should regional biodiversity be monitored?
762      *Environmental and Ecological Statistics* **19**(4), 601-626.
763 BUCKLAND, S. T. & JOHNSTON, A. (2017). Monitoring the biodiversity of regions: Key principles and possible
764      pitfalls. *Biological Conservation* **214**, 23-34.
765 BUSH, A., SOLLMANN, R., WILTING, A., BOHMANN, K., COLE, B., BALZTER, H., MARTIUS, C., ZLINSZKY, A., CALVIGNAC-
766      SPENCER, S., COBBOLD, C. A., DAWSON, T. P., EMERSON, B. C., FERRIER, S., GILBERT, M. T. P., HEROLD, M.,
767      JONES, L., LEENDERTZ, F. H., MATTHEWS, L., MILLINGTON, J. D. A., OLSON, J. R., OVASKAINEN, O., RAFFAELLI,
768      D., REEVE, R., RODEL, M. O., RODGERS, T. W., SNAPE, S., VISSEREN-HAMAKERS, I., VOGLER, A. P., WHITE, P.
769      C. L., WOOSTER, M. J. & YU, D. W. (2017). Connecting Earth observation to high-throughput
770      biodiversity data. *Nature Ecology & Evolution* **1**(7).
771 CALLAGHAN, C. T., POORE, A. G. B., MAJOR, R. E., ROWLEY, J. J. L. & CORNWELL, W. K. (2019). Optimizing future
772      biodiversity sampling by citizen scientists. *Proceedings of the Royal Society B-Biological Sciences*
773      **286**(1912).
774 CALLAGHAN, C. T., THOMPSON, M., WOODS, A., POORE, A. G. B., BOWLER, D. E., SAMONTE, F., ROWLEY, J. J. L.,
775      ROSLAN, N., KINGSFORD, R. T., CORNWELL, W. K. & MAJOR, R. E. (2023). Experimental evidence that
776      behavioral nudges in citizen science projects can improve biodiversity data. *Bioscience* **73**(4),
777      302-313.
778 CALLCUTT, K., CROFT, S. & SMITH, G. C. (2018). Predicting population trends using citizen science data: do
779      subsampling methods produce reliable estimates for mammals? *European Journal of Wildlife*
780      *Research* **64**(3).
781 CARDINALE, B. J., GONZALEZ, A., ALLINGTON, G. R. H. & LOREAU, M. (2018). Is local biodiversity declining or
782      not? A summary of the debate over analysis of species richness time trends. *Biological*
783      *Conservation* **219**, 175-183.
784 CARPENTER, J. & KENWARD, M. (2012). *Multiple imputaion and its application*. Wiley.
785 CAUGHEY, D., BERINSKY, A. J., CHATFIELD, S., HARTMAN, E., SCHICKLER, E. & SEKHON, J. S. (2020). *Target*
786      *estimation and adjustment weighting for survey nonresponse and sampling bias*. Cambridge
787      University Press.
788 CHANDLER, M., SEE, L., COPAS, K., BONDE, A. M. Z., LOPEZ, B. C., DANIELSEN, F., LEGIND, J. K., MASINDE, S., MILLER-
789      RUSHING, A. J., NEWMAN, G., ROSEMARTIN, A. & TURAK, E. (2017). Contribution of citizen science
790      towards international biodiversity monitoring. *Biological Conservation* **213**, 280-294.

791    CHEVALIER, M., ZARZO-ARIAS, A., GUELAT, J., MATEO, R. G. & GUISAN, A. (2022). Accounting for niche
792        truncation to improve spatial and temporal predictions of species distributions. *Frontiers in*
793        *Ecology and Evolution* **10**.
794    CHIFFARD, J., MARCIAU, C., YOCCOZ, N. G., MOUILLOT, F., DUCHATEAU, S., NADEAU, I., FONTANILLES, P. & BESNARD,
795        A. (2020). Adaptive niche-based sampling to improve ability to find rare and elusive species:
796        Simulations and field tests. *Methods in Ecology and Evolution* **11**(8), 899-909.
797    COLLINS, L. M., SCHAFER, J. L. & KAM, C. M. (2001). A comparison of inclusive and restrictive strategies in
798        modern missing data procedures. *Psychological Methods* **6**(4), 330-351.
799    CONN, P. B., THORSON, J. T. & JOHNSON, D. S. (2017). Confronting preferential sampling when analysing
800        population distributions: diagnosis and model-based triage. *Methods in Ecology and Evolution*
801        **8**(11), 1535-1546.
802    COURTER, J. R., JOHNSON, R. J., STUYCK, C. M., LANG, B. A. & KAISER, E. W. (2013). Weekend bias in Citizen
803        Science data reporting: implications for phenology studies. *International Journal of*
804        *Biometeorology* **57**(5), 715-720.
805    CRETOIS, B., SIMMONDS, E. G., LINNELL, J. D. C., VAN MOORTER, B., ROLANDSEN, C. M., SOLBERG, E. J., STRAND, O.,
806        GUNDERSEN, V., ROER, O. & ROD, J. K. (2021). Identifying and correcting spatial bias in opportunistic
807        citizen science data for wild ungulates in Norway. *Ecology and Evolution* **11**(21), 15191-15204.
808    DAKKI, M., ROBIN, G., SUET, M., QNINBA, A., EL AGBANI, M. A., OUASSOU, A., EL HAMOUMI, R., AZAFZAF, H., REBAH,
809        S., FELTRUP-AZAFZAF, C., HAMOUDA, N., IBRAHIM, W. A. L., ASRAN, H. H., ELHADY, A. A., IBRAHIM, H.,
810        ETAYEB, K., BOURAS, E., SAIED, A., GLIDAN, A., HABIB, B. M., SAYOUD, M. S., BENDJEDDA, N., DAMI, L.,
811        DESCHAMPS, C., GAGET, E., MONDAIN-MONVAL, J. Y. & DU RAU, P. D. (2021). Imputation of incomplete
812        large-scale monitoring count data via penalized estimation. *Methods in Ecology and Evolution*
813        **12**(6), 1031-1039.
814    DAMBLY, L. I., JONES, K. E., BOUGHEY, K. L. & ISAAC, N. J. B. (2021). Observer retention, site selection and
815        population dynamics interact to bias abundance trends in bats. *Journal of Applied Ecology* **58**(2),
816        236-247.
817    DARVILL, B., HARRIS, S. J., MARTAY, B., WILSON, M. & GILLINGS, S. (2020). Delivering robust population trends
818        for Scotland's widespread breeding birds. *Scottish Birds* **40**(4), 297-304.
819    DENNIS, E. B., MORGAN, B. J. T., FREEMAN, S. N., BRERETON, T. M. & ROY, D. B. (2016). A Generalized
820        Abundance Index for Seasonal Invertebrates. *Biometrics* **72**(4), 1305-1314.
821    DIEKERT, F., MUNZINGER, S., SCHULEMANN-MAIER, G. & STADTLER, L. (2023). Explicit incentives increase citizen
822        science recordings. *Conservation Letters*.
823    DIGGLE, P. J., MENEZES, R. & SU, T. L. (2010). Geostatistical inference under preferential sampling. *Journal*
824        *of the Royal Statistical Society Series C-Applied Statistics* **59**, 191-232.
825    DORNELAS, M., GOTELLI, N. J., McGILL, B., SHIMADZU, H., MOYES, F., SIEVERS, C. & MAGURRAN, A. E. (2014).
826        Assemblage Time Series Reveal Biodiversity Change but Not Systematic Loss. *Science* **344**(6181),
827        296-299.
828    EVANS, D. M. & DAY, K. R. (2002). Hunting disturbance on a large shallow lake: the effectiveness of
829        waterfowl refuges. *Ibis* **144**(1), 2-8.
830    FINK, D., AUER, T., JOHNSTON, A., RUIZ-GUTIERREZ, V., HOCHACHKA, W. M. & KELLING, S. (2020). Modeling avian
831        full annual cycle distribution and population trends with citizen science data. *Ecological*
832        *Applications* **30**(3).
833    FITHIAN, W., ELITH, J., HASTIE, T. & KEITH, D. A. (2015). Bias correction in species distribution models: pooling
834        survey and collection data for multiple species. *Methods in Ecology and Evolution* **6**(4), 424-438.
835    FORISTER, M. L., BLACK, S. H., ELPHICK, C. S., GRAMES, E. M., HALSCH, C. A., SCHULTZ, C. B. & WAGNER, D. L. (2023).
836        Missing the bigger picture: Why insect monitoring programs are limited in their ability to
837        document the effects of habitat loss. *Conservation Letters* **16**(3).

838    FRAIR, J. L., NIELSEN, S. E., MERRILL, E. H., LELE, S. R., BOYCE, M. S., MUNRO, R. H. M., STENHOUSE, G. B. & BEYER,
839        H. L. (2004). Removing GPS collar bias in habitat selection studies. *Journal of Applied Ecology*
840        **41**(2), 201-212.
841    FRAISL, D., CAMPBELL, J., SEE, L., WEHN, U., WARDLAW, J., GOLD, M., MOORTHY, I., ARIAS, R., PIERA, J., OLIVER, J. L.,
842        MASO, J., PENKER, M. & FRITZ, S. (2020). Mapping citizen science contributions to the UN
843        sustainable development goals. *Sustainability Science* **15**(6), 1735-1751.
844    FRETWELL, P. T., SCOFIELD, P. & PHILLIPS, R. A. (2017). Using super-high resolution satellite imagery to census
845        threatened albatrosses. *Ibis* **159**(3), 481-490.
846    GARCIA-ROSELLO, E., GUISANDE, C., MANJARRES-HERNANDEZ, A., GONZALEZ-DACOSTA, J., HEINE, J., PELAYO-VILLAMIL,
847        P., GONZALEZ-VILAS, L., VARI, R. P., VAAMONDE, A., GRANADO-LORENCIO, C. & LOBO, J. M. (2015). Can we
848        derive macroecological patterns from primary Global Biodiversity Information Facility data?
849        *Global Ecology and Biogeography* **24**(3), 335-347.
850    GAUL, W., SADYKOVA, D., WHITE, H. J., LEON-SANCHEZ, L., CAPLAT, P., EMMERSON, M. C. & YEARSLEY, J. M. (2022).
851        Modelling the distribution of rare invertebrates by correcting class imbalance and spatial bias.
852        *Diversity and Distributions* **28**(10), 2171-2186.
853    GELDMANN, J., HEILMANN-CLAUSEN, J., HOLM, T. E., LEVINSKY, I., MARKUSSEN, B., OLSEN, K., RAHBEK, C. & TOTTRUP,
854        A. P. (2016). What determines spatial bias in citizen science? Exploring four recording schemes
855        with different proficiency requirements. *Diversity and Distributions* **22**(11), 1139-1149.
856    GELMAN, A. (2007). Struggles with survey weighting and regression modeling. *Statistical Science* **22**(2),
857        153-164.
858    GONZALEZ, A., CARDINALE, B. J., ALLINGTON, G. R. H., BYRNES, J., ENDSLEY, K. A., BROWN, D. G., HOOPER, D. U.,
859        ISBELL, F., O'CONNOR, M. I. & LOREAU, M. (2016). Estimating local biodiversity change: a critique of
860        papers claiming no net loss of local diversity. *Ecology* **97**(8), 1949-1960.
861    GRATTAROLA, F., BOWLER, D. E. & KEIL, P. (2023). Integrating presence-only and presence-absence data to
862        model changes in species geographic ranges: An example in the Neotropics. *Journal of*
863        *Biogeography*.
864    GREGORY, R. D., VAN STRIEN, A., VORISEK, P., MEYLING, A. W. G., NOBLE, D. G., FOPPEN, R. P. B. & GIBBONS, D. W.
865        (2005). Developing indicators for European birds. *Philosophical Transactions of the Royal Society*
866        *B-Biological Sciences* **360**(1454), 269-288.
867    GUO, Q. F., CHEN, A. P., CROCKETT, E. T. H., ATKINS, J. W., CHEN, X. W. & FEI, S. L. (2023). Integrating gradient
868        with scale in ecological and evolutionary studies. *Ecology* **104**(4).
869    HEFLEY, T. J., TYRE, A. J., BAASCH, D. M. & BLANKENSHIP, E. E. (2013). Nondetection sampling bias in marked
870        presence-only data. *Ecology and Evolution* **3**(16), 5225-5236.
871    HERRERA, C. M. (2019). Complex long-term dynamics of pollinator abundance in undisturbed
872        Mediterranean montane habitats over two decades. *Ecological Monographs* **89**(1).
873    HERTZOG, L. R., BESNARD, A. & JAY-ROBERT, P. (2014). Field validation shows bias-corrected pseudo-absence
874        selection is the best method for predictive species-distribution modelling. *Diversity and*
875        *Distributions* **20**(12), 1403-1413.
876    HOF, A. R. & BRIGHT, P. W. (2016). Quantifying the long-term decline of the West European hedgehog in
877        England by subsampling citizen-science datasets. *European Journal of Wildlife Research* **62**(4),
878        407-413.
879    HOSSIE, T. J., GOBIN, J. & MURRAY, D. L. (2021). Confronting Missing Ecological Data in the Age of Pandemic
880        Lockdown. *Frontiers in Ecology and Evolution* **9**.
881    ISAAC, N. J. B. & POCOCK, M. J. O. (2015). Bias and information in biological records. *Biological Journal of*
882        *the Linnean Society* **115**(3), 522-531.
883    JETZ, W., MCGEOCH, M. A., GURALNICK, R., FERRIER, S., BECK, J., COSTELLO, M., FERNANDEZ, M., GELLER, G. N., KEIL,
884        P., MEROW, C., MEYER, C., MULLER-KARGER, F. E., PEREIRA, H. M., REGAN, E. C., SCHMELLER, D. S. &

885     TURAK, E. (2019). Essential biodiversity variables for mapping and monitoring species
886          populations. *Nature Ecology & Evolution* **3**(4), 539-551.
887 JOHNSON, T. F., ISAAC, N. J. B., PAVIOLO, A. & GONZÁLEZ-SUÁREZ, M. (2023). Socioeconomic factors predict
888          population changes of large carnivores better than climate change or habitat loss. . 2023 Jan
889          24;14(1):74. doi: . PMID: 36693827; PMCID: PMC9873912. *Nat Commun* **14**(1), 74.
890 JOHNSTON, A., HOCHACHKA, W. M., STRIMAS-MACKEY, M. E., GUTIERREZ, V. R., ROBINSON, O. J., MILLER, E. T., AUER,
891          T., KELLING, S. T. & FINK, D. (2021). Analytical guidelines to increase the value of community
892          science data: An example using eBird data to estimate species distributions. *Diversity and*
893          *Distributions* **27**(7), 1265-1277.
894 JOHNSTON, A., MORAN, N., MUSGROVE, A., FINK, D. & BAILLIE, S. R. (2020). Estimating species distributions
895          from spatially biased citizen science data. *Ecological Modelling* **422**.
896 KISSLING, W. D., AHUMADA, J. A., BOWSER, A., FERNANDEZ, M., FERNANDEZ, N., GARCIA, E. A., GURALNICK, R. P.,
897          ISAAC, N. J. B., KELLING, S., LOS, W., MCRAE, L., MIHOUB, J. B., OBST, M., SANTAMARIA, M., SKIDMORE, A.
898          K., WILLIAMS, K. J., AGOSTI, D., AMARILES, D., ARVANITIDIS, C., BASTIN, L., DE LEO, F., EGLOFF, W., ELITH, J.,
899          HOBERN, D., MARTIN, D., PEREIRA, H. M., PESOLE, G., PETERSEIL, J., SAARENMAA, H., SCHIGEL, D.,
900          SCHMELLER, D. S., SEGATA, N., TURAK, E., UHLIR, P. F., WEE, B. & HARDISTY, A. R. (2018). Building
901          essential biodiversity variables (EBVs) of species distribution and abundance at a global scale.
902          *Biological Reviews* **93**(1), 600-625.
903 KREFT, H. & JETZ, W. (2007). Global patterns and determinants of vascular plant diversity. *Proceedings of*
904          *the National Academy of Sciences of the United States of America* **104**(14), 5925-5930.
905 LA SORTE, F. A. & SOMVEILLE, M. (2020). Survey completeness of a global citizen-science database of bird
906          occurrence. *Ecography* **43**(1), 34-43.
907 LANEY, J. A., HALLMAN, T. A., CURTIS, J. R. & ROBINSON, W. D. (2021). The influence of rare birds on observer
908          effort and subsequent rarity discovery in the American birdwatching community. *Peerj* **9**.
909 LAXTON, M. R., DE RIVERA, O. R., SORIANO-REDONDO, A. & ILLIAN, J. B. (2023). Balancing structural complexity
910          with ecological insight in Spatio-temporal species distribution models. *Methods in Ecology and*
911          *Evolution* **14**(1), 162-172.
912 LEE, K. J., TILLING, K. M., CORNISH, R. P., LITTLE, R. J. A., BELL, M. L., GOETGHEBEUR, E., HOGAN, J. W., CARPENTER, J.
913          R. & INITIATIVE, S. (2021). Framework for the treatment and reporting of missing data in
914          observational studies: The Treatment And Reporting of Missing data in Observational Studies
915          framework. *Journal of Clinical Epidemiology* **134**, 79-88.
916 LEHIKOINEN, A., FOPPEN, R. P. B., HELDBJERG, H., LINDSTROM, A., VAN MANEN, W., PIIRAINEN, S., VAN TURNHOUT, C.
917          A. M. & BUTCHART, S. H. M. (2016). Large-scale climatic drivers of regional winter bird population
918          trends. *Diversity and Distributions* **22**(11), 1163-1173.
919 LEURENT, B., GOMES, M., FARIA, R., MORRIS, S., GRIEVE, R. & CARPENTER, J. R. (2018). Sensitivity Analysis for
920          Not-at-Random Missing Data in Trial-Based Cost-Effectiveness Analysis: A Tutorial.
921          *Pharmacoeconomics* **36**(8), 889-901.
922 LI, L. L., SHEN, C. Y., LI, X. C. & ROBINS, J. M. (2013). On weighting approaches for missing data. *Statistical*
923          *Methods in Medical Research* **22**(1), 14-30.
924 LIN, Y. P., YEH, M. S., DENG, D. P. & WANG, Y. C. (2008). Geostatistical approaches and optimal additional
925          sampling schemes for spatial patterns and future sampling of bird diversity. *Global Ecology and*
926          *Biogeography* **17**(2), 175-188.
927 LITTLE, R. J., CARPENTER, J. R., LEE, K. J. & INITIATIVE, S. (2022). A Comparison of Three Popular Methods for
928          Handling Missing Data: Complete-Case Analysis, Inverse Probability Weighting, and Multiple
929          Imputation. *Sociological Methods & Research*.
930 LITTLE, R. J. & RUBIN, D. B. (2019). *Statistical Analysis with Missing Data (3rd Edition)*. Wiley.
931 LITTLE, R. J. & VARTIVARIAN, S. (2005). Does Weighting for Nonresponse Increase the Variance of Survey
932          Means? . *Survey Methodology* **31**(2), 161-168.

933    LITTLE, R. J. A. (1988). A Test of Missing Completely at Random for Multivariate Data with Missing Values.
934         *Journal of the American Statistical Association* **83**(404), 1198-1202.
935    LITTLE, R. J. A. (1995). Modeling the Drop-out Mechanism in Repeated-Measures Studies. *Journal of the*
936         *American Statistical Association* **90**(431), 1112-1121.
937    LITTLE, T. D. & RHEMTULLA, M. (2013). Planned Missing Data Designs for Developmental Researchers. *Child*
938         *Development Perspectives* **7**(4), 199-204.
939    LOPUCKI, R., KIERSZTYN, A., PITUCHA, G. & KITOWSKI, I. (2022). Handling missing data in ecological studies:
940         Ignoring gaps in the dataset can distort the inference. *Ecological Modelling* **468**.
941    MANDEVILLE, C. P., NILSEN, E. B. & FINSTAD, A. G. (2022). Spatial distribution of biodiversity citizen science in
942         a natural area depends on area accessibility and differs from other recreational area use.
943         *Ecological Solutions and Evidence* **3**(4).
944    MARSH, D. M. & COSENTINO, B. J. (2019). Causes and consequences of non-random drop-outs for citizen
945         science projects: lessons from the North American amphibian monitoring program. *Freshwater*
946         *Science* **38**(2), 292-302.
947    MATUTINI, F., BAUDRY, J., PAIN, G., SINEAU, M. & PITHON, J. (2021). How citizen science could improve species
948         distribution models and their independent assessment. *Ecology and Evolution* **11**(7), 3028-3039.
949    MCCLURE, C. J. W. & ROLEK, B. W. (2023). Pitfalls arising from site selection bias in population monitoring
950         defy simple heuristics. *Methods in Ecology and Evolution* **14**(6), 1489-1499.
951    MENG, X. L. (2018). Statistical Paradises and Paradoxes in Big Data (I): Law of Large Populations, Big Data
952         Paradox, and the 2016 Us Presidential Election. *Annals of Applied Statistics* **12**(2), 685-726.
953    MENG, X. L. (2022). Comments on "Statistical inference with non-probability survey samples" -
954         Miniaturizing data defect correlation: A versatile strategy for handling non-probability samples.
955         *Survey Methodology* **48**(2), 339-360.
956    MOHAN, K. & PEARL, J. (2021). Graphical Models for Processing Missing Data. *Journal of the American*
957         *Statistical Association* **116**(534), 1023-1037.
958    NAKAGAWA, S. & FRECKLETON, R. P. (2008). Missing inaction: the dangers of ignoring missing data. *Trends in*
959         *Ecology & Evolution* **23**(11), 592-596.
960    NIELSEN, S. E., HAUGHLAND, D. L., BAYNE, E. & SCHIECK, J. (2009). Capacity of large-scale, long-term
961         biodiversity monitoring programmes to detect trends in species prevalence. *Biodiversity and*
962         *Conservation* **18**(11), 2961-2978.
963    NOBLE, D. W. A. & NAKAGAWA, S. (2021). Planned missing data designs and methods: Options for
964         strengthening inference, increasing research efficiency and improving animal welfare in
965         ecological and evolutionary research. *Evolutionary Applications* **14**(8), 1958-1968.
966    NUNEZ-PENICHET, C., COBOS, M. E., SOBERON, J., GUETA, T., BARVE, N., BARVE, V., NAVARRO-SIGUENZA, A. G. &
967         PETERSON, A. T. (2022). Selection of sampling sites for biodiversity inventory: Effects of
968         environmental and geographical considerations. *Methods in Ecology and Evolution* **13**(7), 1595-
969         1607.
970    OUTHWAITE, C. L., POWNEY, G. D., AUGUST, T. A., CHANDLER, R. E., RORKE, S., PESCOTT, O. L., HARVEY, M., ROY, H.
971         E., FOX, R., ROY, D. B., ALEXANDER, K., BALL, S., BANTOCK, T., BARBER, T., BECKMANN, B. C., COOK, T.,
972         FLANAGAN, J., FOWLES, A., HAMMON, P., HARVEY, P., HEPPER, D., HUBBLE, D., KRAMER, J., LEE, P.,
973         MACADAM, C., MORRIS, R., NORRIS, A., PALMER, S., PLANT, C. W., SIMKIN, J., STUBBS, A., SUTTON, P.,
974         TELFER, M., WALLACE, I. & ISAAC, N. J. B. (2019). Annual estimates of occupancy for bryophytes,
975         lichens and invertebrates in the UK, 1970-2015. *Scientific Data* **6**.
976    PENNINO, M. G., PARADINAS, I., ILLIAN, J. B., MUNOZ, F., BELLIDO, J. M., LOPEZ-QUILEZ, A. & CONESA, D. (2019).
977         Accounting for preferential sampling in species distribution models. *Ecology and Evolution* **9**(1),
978         653-663.
979    POCOCK, M. J., TWEDDLE, J. C., SAVAGE, J., ROBINSON, L. D. & ROY, H. E. (2017). The diversity and evolution of
980         ecological and environmental citizen science. *Plos One* **12**(4).

981    RAGHUNATHAN, T. E. (2004). What do we do with missing data? Some options for analysis of incomplete
982            data. *Annual Review of Public Health* **25**, 99-117.
983    RAPACCIUOLO, G., YOUNG, A. & JOHNSON, R. (2021). Deriving indicators of biodiversity change from
984            unstructured community-contributed data. *Oikos* **130**(8), 1225-1239.
985    REHFISCH, M. M., AUSTIN, G. E., ARMITAGE, M. J. S., ATKINSON, P. W., HOLLOWAY, S. J., MUSGROVE, A. J. & POLLITT,
986            M. S. (2003). Numbers of wintering waterbirds in Great Britain and the Isle of Man (1994/1995-
987            1998/1999): II. Coastal waders (Charadrii). *Biological Conservation* **112**(3), 329-341.
988    ROBINSON, O. J., RUIZ-GUTIERREZ, V., REYNOLDS, M. D., GOLET, G. H., STRIMAS-MACKEY, M. & FINK, D. (2020).
989            Integrating citizen science data with expert surveys increases accuracy and spatial extent of
990            species distribution models. *Diversity and Distributions* **26**(8), 976-986.
991    RUBIN, D. B. (1976). Inference and Missing Data. *Biometrika* **63**(3), 581-590.
992    SCHMUCKI, R., PE'ER, G., ROY, D. B., STEFANESCU, C., VAN SWAAY, C. A. M., OLIVER, T. H., KUUSSAARI, M., VAN
993            STRIEN, A. J., RIES, L., SETTELE, J., MUSCHE, M., CARNICER, J., SCHWEIGER, O., BRERETON, T. M., HARPKE, A.,
994            HELIOLA, J., KUHN, E. & JULLIARD, R. (2016). A regionally informed abundance index for supporting
995            integrative analyses across butterfly monitoring schemes. *Journal of Applied Ecology* **53**(2), 501-
996            510.
997    SEAMAN, S. R. & WHITE, I. R. (2013). Review of inverse probability weighting for dealing with missing data.
998            *Statistical Methods in Medical Research* **22**(3), 278-295.
999    SMITH, T. M. F. (1976). The Foundations of Survey Sampling *Journal of the Royal Statistical Societiy Series*
1000          *A* **139**, 183-204.
1001    SPAKE, R., BOWLER, D. E., CALLAGHAN, C. T., BLOWES, S. A., DONCASTER, C. P., ANTAO, L. H., NAKAGAWA, S.,
1002            MCELREATH, R. & CHASE, J. M. (2023). Understanding 'it depends' in ecology: a guide to
1003            hypothesising, visualising and interpreting statistical interactions. *Biological Reviews* **98**(4), 983-
1004            1002.
1005    SPAKE, R., O'DEA, R. E., NAKAGAWA, S., DONCASTER, C. P., RYO, M., CALLAGHAN, C. T. & BULLOCK, J. M. (2022).
1006            Improving quantitative synthesis to achieve generality in ecology. *Nature Ecology & Evolution*
1007            **6**(12), 1818-1828.
1008    SPECHT, H. M., REICH, H. T., IANNARILLI, F., EDWARDS, M. R., STAPLETON, S. P., WEEGMAN, M. D., JOHNSON, M. K.,
1009            YOHANNES, B. J. & ARNOLD, T. W. (2017). Occupancy surveys with conditional replicates: An
1010            alternative sampling design for rare species. *Methods in Ecology and Evolution* **8**(12), 1725-1734.
1011    STEEN, V. A., ELPHICK, C. S. & TINGLEY, M. W. (2019). An evaluation of stringent filtering to improve species
1012            distribution models from citizen science data. *Diversity and Distributions* **25**(12), 1857-1869.
1013    STEEN, V. A., TINGLEY, M. W., PATON, P. W. C. & ELPHICK, C. S. (2021). Spatial thinning and class balancing:
1014           Key choices lead to variation in the performance of species distribution models with citizen
1015           science data. *Methods in Ecology and Evolution* **12**(2), 216-226.
1016    SULLIVAN, B. L., PHILLIPS, T., DAYER, A. A., WOOD, C. L., FARNSWORTH, A., ILIFF, M. J., DAVIES, I. J., WIGGINS, A.,
1017            FINK, D., HOCHACHKA, W. M., RODEWALD, A. D., ROSENBERG, K. V., BONNEY, R. & KELLING, S. (2017).
1018            Using open access observational data for conservation action: A case study for birds. *Biological*
1019            *Conservation* **208**, 5-14.
1020    TCHETGEN, E. J. T. & WIRTH, K. E. (2017). A general instrumental variable framework for regression analysis
1021            with outcome missing not at random. *Biometrics* **73**(4), 1123-1131.
1022    TER BRAAK, C. J. F., VAN STRIEN, A. J., MEIJER, R. & VERSTRAEL, T. J. (1992). *Analysis of monitoring data with*
1023          *many missing values: which method?*
1024    THOEMMES, F. & ROSE, N. (2014). A Cautious Note on Auxiliary Variables That Can Increase Bias in Missing
1025            Data Problems. *Multivariate Behavioral Research* **49**(5), 443-459.
1026    TROUDET, J., GRANDCOLAS, P., BLIN, A., VIGNES-LEBBE, R. & LEGENDRE, F. (2017). Taxonomic bias in biodiversity
1027            data and societal preferences. *Scientific Reports* **7**.

1028    Tulloch, A. I. T., Mustin, K., Possingham, H. P., Szabo, J. K. & Wilson, K. A. (2013). To boldly go where no
1029        volunteer has gone before: predicting volunteer activity to prioritize surveys at the landscape
1030        scale. *Diversity and Distributions* **19**(4), 465-480.
1031    Underhill, L. G. & Prysjones, R. P. (1994). Index Numbers for Waterbird Populations .1. Review and
1032        Methodology. *Journal of Applied Ecology* **31**(3), 463-480.
1033    Valdez, J. W., Callaghan, C. T., Junker, J., Purvis, A., Hill, S. L. L. & Pereira, H. M. (2023). The
1034        undetectability of global biodiversity trends using local species richness. *Ecography* **2023**(3).
1035    Valliant, R. (2020). Comparing Alternatives for Estimation from Nonprobability Samples. *Journal of
1036        Survey Statistics and Methodology* **8**(2), 231-263.
1037    van den Brakel, J. A. & Bethlehem, J. (2008). *Model-based estimation for official statistics*. Statistics
1038        Netherlands, Voorburg/Heerlen.
1039    van Swaay, C. A. M., Nowicki, P., Settele, J. & van Strien, A. J. (2008). Butterfly monitoring in Europe:
1040        methods, applications and perspectives. *Biodiversity and Conservation* **17**(14), 3455-3469.
1041    Ver Hoef, J. M., Johnson, D., Angliss, R. & Higham, M. (2021). Species density models from opportunistic
1042        citizen science data. *Methods in Ecology and Evolution* **12**(10), 1911-1925.
1043    Wu, C. B. (2022). Statistical inference with non-probability survey samples. *Survey Methodology* **48**(2),
1044        283-311.
1045    Zbinden, N., Kery, M., Hafliger, G., Schmid, H. & Keller, V. (2014). A resampling-based method for effort
1046        correction in abundance trend analyses from opportunistic biological records. *Bird Study* **61**(4),
1047        506-517.
1048    Zhang, W. Y., Sheldon, B., Grenyer, R. & Gaston, K. J. (2021). Habitat change and biased sampling
1049        influence estimation of diversity trends. *Current Biology* **31**(16), 3656-+.
1050    Zimney, A. & Smart, T. (2022). Effects of incomplete sampling and standardization on indices of
1051        abundance from a fishery- independent trawl survey off the Atlantic coast of the southeastern
1052        United States. *Fishery Bulletin* **120**(3-4), 252-267.

1053

# Supporting Information

**Table S1 Selected R tools that can help with missing data problems and their potential application for use in biodiversity research.**

| R packages | Applications | Useful functions |
|---|---|---|
| *Exploring missing data* | | |
| naniar | visualizing/exploring the missing data pattern | *mcar_test* – Little's missing completely at random (MCAR) test<br>*vis_miss* – plot the missing data for all variables |
| occAssess | measure of the potential for bias in taxonomic, temporal, spatial, and environmental dimensions | *assessEnvBias* - assess whether data are sampled from a representative portion of environmental space in the spatial domain of interest<br>*assessSpatialBias* – assess whether data resemble a random distribution in the geographic space of interest for inference<br>*assessSpatialCov* – assess whether a representative portion of the spatial domain of interest has been sampled and whether the same portion of geographic space has been sampled over time |
| sampbias | a Bayesian approach to estimate how sampling rates vary as a function of proximity to one or multiple bias factors | *calculate_bias* - calculating the bias effect of sampling bias due to geographic structures, such as the vicinity to cities, airports, rivers and roads |
| *Subsampling* | | |
| base | Base R functions | *sample* - sample data with predefined inclusion probabilities specified with the prob argument |
| sampling | draw random samples using different sampling schemes | *balancedcluster* – selects a balanced cluster sample according to defined auxiliary variables<br>*strata* - stratified sampling with unequal probabilities. |

| spatialEco | spatial data manipulation and modelling | *stratified.random* - creates a stratified random sample of an sp class object <br> *stratified.distance* - draws a minimum, and optional maximum constrained, distance sub-sampling |
|---|---|---|
| spThin | Spatial thinning of species occurence records | *thin* - returns a dataset with the maximum number of records for a given thinning distance |
| terra | spatial data manipulation and processing | *spatSample* – sample a SpatRaster, SpatVector or SpatExtent objcy |
| *Imputation* | | |
| agTrend | modelling regional trends with missing data | *mcmc.aggregate* - a zero-inflated, nonparameteric model with a definable observation model, augmenting missing values before calculating regional abundances |
| INLA/<br>inlabru | fitting Bayesian models, especially useful for spatial models via its spatial mesh | *inla/bru* - fit a Bayesian model using Integrated Nested Laplace approximation <br> *predict* – draw predictions from the fitted model, where the prediction data frame can be a SpatialPointsDataFrame object |
| LORI | imputation of missing count data | *lori* – impute missing count data using a large covariate set, including interactions, with a LASSO penalty |
| mice | multiple imputation by chained equations | *mice* – multiple imputation method that will generate plausible values for any missing data – in the response and in any covariates |
| Rjags<br>JAGS<br>nimble | fitting Bayesian models allowing for missing values in the response | *Jags/runMCMC* - fitting Bayesian models allowing for imputation of missing values in the response during model fitting (options available for missing values in covariates too) |
| rtrim | functions to calculate annual indices and trends of abundances | fit a GLM imputing missing values based on mean site and year effects, with optional covariates |
| *Weighting* | | |

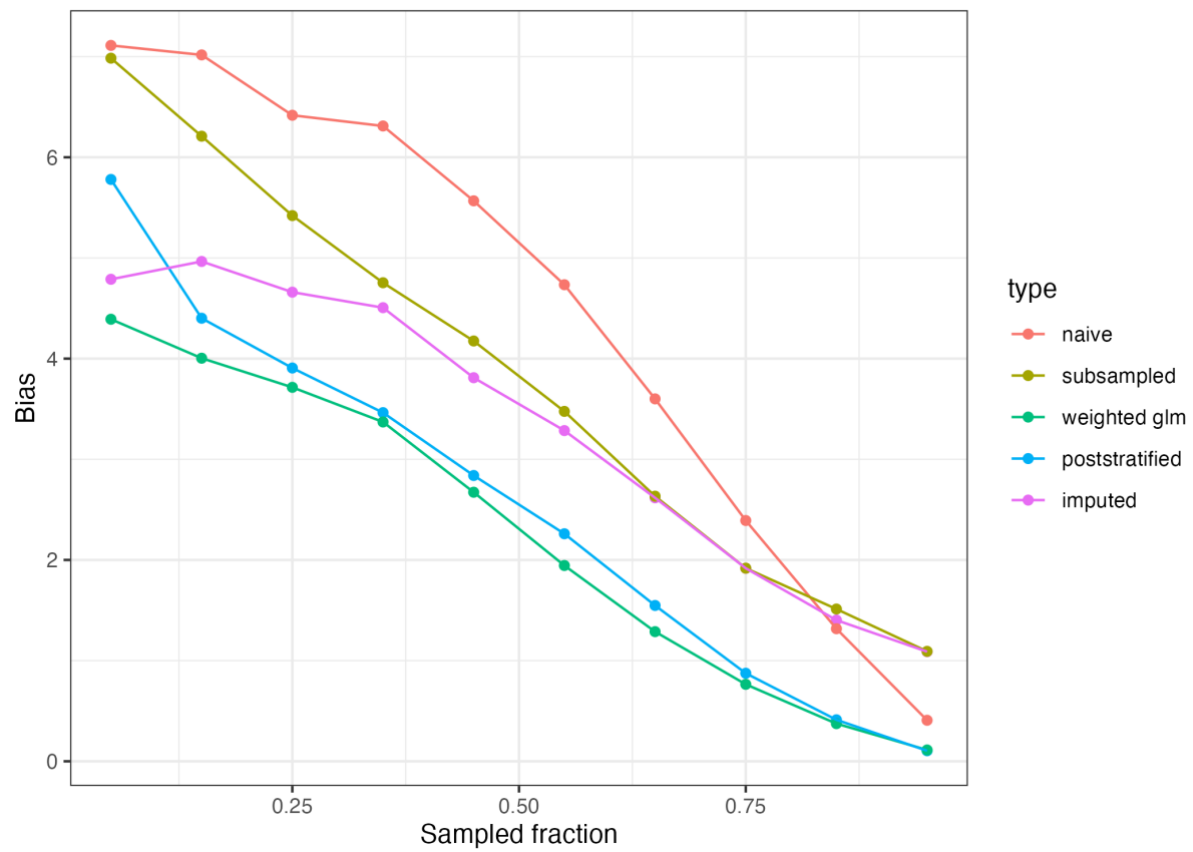| survey srvyr | range of functions for analysis of data from complex surveys, including fitting models with weights | *Svyglm* – generalized linear models with survey weights<br>*postStratify* – function for post-stratification to match the joint distribution of the variables of the population |
|---|---|---|
| svrep | Analysis of replicate/boostrapped survey weights | *svyby_repwts* – compare estimates from different sets of weights |
| twang | functions to estimate propensity scores and weights | *ps* - gradient boosted trees to predict non-response from covariates<br>*bal.table* – compare covariate values between sample and population |

Fig. S1 The ability of missing data solutions to adjust for bias in biodiversity data.

We assumed a landscape of 400 cells and that a covariate affected both species abundance and the likelihood of a cell being sampled. We vary the fraction of the cells that were sampled. In contrast to Fig 5A (main text), we assumed that the species abundance was affected by an additional covariate that did not affect sampling; this variable was not included in any of the analysis. The models to estimate the parameter of interest (mean abundance) were: naive (no correction); subsampled (cells were subsampled along the covariate gradient to reduce the sampling bias), weighted (two methods: weighted glm using the svyglm function, and weighted by poststratification, using postStratify, both in the survey package) and imputed (using JAGS to impute NAs in the response). Points show the mean bias (difference between model prediction and truth) across 100 independent runs.