

1 **Descriptive inference using large, unrepresentative nonprobability samples: An**
2 **introduction for ecologists**

3 ^{*1}Robin J. Boyd, ²Gavin B. Stewart and ¹Oliver L. Pescott

4 ¹UK Centre for Ecology & Hydrology, Benson Lane, Wallingford, OX108BB

5 ² Evidence Synthesis Lab, School of Natural and Environmental Science, University of
6 Newcastle, Newcastle-upon-Tyne, NE1 7RU

7 *corresponding author email: robboy@ceh.ac.uk

8 **Abstract**

9 Biodiversity monitoring usually involves drawing inferences about some variable of interest
10 across a defined landscape from observations made at a sample of locations within that
11 landscape. If the variable of interest differs between sampled and non-sampled locations, and
12 no mitigating action is taken, then the sample is unrepresentative and inferences drawn from
13 it will be biased. It is possible to adjust unrepresentative samples so that they more closely
14 resemble the wider landscape in terms of “auxiliary variables”. A good auxiliary variable is a
15 common cause of sample inclusion and the variable of interest, and if it explains an
16 appreciable portion of the variance in both, then inferences drawn from the adjusted sample
17 will be closer to the truth. We applied six types of survey sample adjustment—subsampling,
18 quasi-randomisation, poststratification, superpopulation modelling, a “doubly robust”
19 procedure, and multilevel regression and poststratification—to a simple two-part biodiversity
20 monitoring problem. The first part was to estimate mean occupancy of the plant *Calluna*
21 *vulgaris* in Great Britain in two time-periods (1987-1999 and 2010-2019); the second was to
22 estimate the difference between the two (i.e. the trend). We estimated the means and trend
23 using large, but (originally) unrepresentative, samples from a citizen science dataset.
24 Compared to the unadjusted estimates, the means and trends estimated using most adjustment
25 methods were more accurate, although standard uncertainty intervals generally did not cover
26 the true values. Completely unbiased inference is not possible from an unrepresentative
27 sample without knowing and having data on all relevant auxiliary variables. Adjustments can
28 reduce the bias if auxiliary variables are available and selected carefully, but the potential for
29 residual bias should be acknowledged and reported.

30 **Introduction**

31 As the data revolution gathers pace, it is not surprising to see “big data” being used to
32 monitor biodiversity. Examples include observations submitted to mobile phone apps by
33 amateur naturalists (Johnston et al., 2022) and digitised specimens from museums and
34 herbaria (Nelson & Ellis, 2019). Such data become bigger still when combined in data
35 aggregators such as the Global Biodiversity Information Facility (GBIF;
36 <https://www.gbif.org/>) or metadatabases such as PREDICTS (Hudson et al., 2014).
37 Unfortunately, quantity of data does not necessarily imply quality of insight.

38 Monitoring biodiversity is typically a matter of descriptive statistical inference. It is
39 inferential in that the goal is to infer something about a target population from a sample of
40 that population (Boyd, Powney, et al., 2023). The population might comprise, say, all areal
41 units across some landscape (“sites”), in which case the sample would be a subset of those
42 sites. The inference is descriptive in that the aim is to describe (rather than explain) a variable
43 of interest in the population. A common example is the proportion of sites occupied by some
44 species (Bowler et al., 2021; Outhwaite et al., 2020; Powney et al., 2019; Stroh et al., 2023;
45 van Strien & van Grunsven, 2023), but there are many others.

46 Of more importance than the size of a sample for descriptive inference is whether it is
47 representative of the population about which inferences are to be drawn (Meng, 2018). In a
48 representative sample, the distribution of the variable of interest is similar to its distribution
49 in the population (Bethlehem et al., 2008). An equivalent definition is that there is little to no
50 correlation between inclusion in the sample and the variable of interest—the “data defect
51 correlation”, or ddc (Meng, 2018). Intuitively, statistics derived from a representative sample,
52 such as means and proportions, will be similar to their population equivalents.

53 Unfortunately, ddc’s are likely to be appreciable in big biodiversity datasets. For one,
54 naturalists preferentially visit and collect data at sites where they are likely to see species that
55 interest them (Bowler et al., 2022; Forister et al., 2023). Where those species’ abundances or
56 distributions are the variables of analytic interest, preferential sampling naturally results in a
57 positive ddc (McClure & Rolek, 2023). On the other hand, naturalists might be constrained to
58 visiting and collecting data in, say, built up areas, which are easier to access than remote
59 locations (Geldmann et al., 2016; Hughes et al., 2020; Mandeville et al., 2022). Built-up areas
60 generally have low quality habitat, meaning that species are less likely to occupy them in
61 large numbers and that the ddc might be negative.

62 Inferences from unrepresentative samples, with appreciable ddc’s, are likely to be misleading.
63 Imagine a researcher who wants to estimate the average abundance of some species across a
64 landscape. An obvious (but naïve) approach would be to calculate its mean abundance across
65 sampled sites and assume that this is similar to its average abundance across the wider
66 landscape. However, if the locations at which the species is most abundant were
67 preferentially sampled, then the sample-based estimate of its mean abundance will be
68 upwardly biased. To use the analogy of Forister et al. (2023), sampled locations would be life
69 rafts; non-sampled locations would be the sinking ship.

70 It is simple to counteract the biasing effect of the ddc if the probability that each site was
71 included in the sample is known; that is, if a *probability sample* is available. In this case,
72 more weight can be placed on the data from sites that were less likely to be included. The
73 effect of this type of weighting is easiest to explain heuristically: the sample is augmented
74 with “copies” of the data from sites that were less likely to be sampled, effectively bringing
75 sample inclusion probabilities across sites to parity. Two variables cannot be correlated if one
76 of them is constant, which means there can be no correlation between the weighted sample
77 inclusion probabilities and the variable of interest across sites. It follows that the ddc, which
78 is the correlation between actual (weighted) sample inclusion and the variable of interest, is
79 zero in expectation (Meng, 2022), and the sample can be considered representative (Lohr,
80 2022). Weighting of this type is known as “design-based” inference, because the inclusion
81 probabilities are a feature of the sampling design.

82 Design-based inference is not applicable for the types of big biodiversity datasets we consider
83 here, because they were not collected according to a probabilistic sampling design. We do not
84 know the probabilities that sites were visited by the collectors of specimens now held in
85 museums and on GBIF. Nor do we know the probabilities that citizen scientists visited and
86 collected data at each site across most landscapes. Matters are simpler when using data from
87 structured monitoring schemes, which often aim for a probability sample (e.g PoMS).
88 However, incomplete uptake of sites that were selected for inclusion (Pescott et al., 2015,
89 2019) means that, in practice, these samples too are non-probabilistic. [Incomplete uptake in
90 biodiversity monitoring is analogous to the issue of non-response in survey sampling (e.g.
91 Bethlehem et al., 2008).] Where sample inclusion probabilities are not known, an alternative
92 to design-based inference is needed.

93 Most approaches to inference from nonprobability samples involve *estimating* the inclusion
94 probabilities. A relatively simple example is poststratification, where the observations (for
95 each site) are split into strata based on covariates, and sites in strata that are underrepresented
96 in the population (based on the population totals of the covariates) are given more weight
97 (Valliant et al., 2018). Using covariates to estimate sample inclusion probabilities is
98 equivalent to adjusting the samples in such a way that the distributions of those covariates in
99 the sample more closely resemble their distributions in the population (i.e. across all sites in
100 the wider landscape). If the covariates affect *both* the variable of interest and sample
101 inclusion, then inferences drawn from the adjusted sample will be closer to the truth than
102 those from the original (naïve) sample. In the context of inference from nonprobability
103 samples, covariates affecting both sample inclusion and the variable of interest, which are not
104 of direct analytic interest themselves, are known as “auxiliary variables” (Thoemmes &
105 Mohan, 2015; Thoemmes & Rose, 2014).

106 Before going further, it is important to note that most approaches to inference from
107 nonprobability samples rest on the bold assumption that the variable of interest is
108 independent of sample inclusion after accounting for the auxiliary variables (Bailey, 2022);
109 that is, non-sampled sites are “Missing At Random” (MAR; Rubin, 1976). If the MAR
110 assumption holds, then unbiased inference is possible. In reality, the MAR assumption is
111 likely to be violated, because data are not available on all relevant auxiliary variables, so the
112 best we can hope for is a reduction in bias relative to naïve inferences drawn from the
113 unadjusted sample.

114 Use of sample adjustments in biodiversity monitoring is variable. It is common for
115 monitoring schemes to weight samples in such a way that the relative frequencies of habitats
116 or geographic areas in the sample are similar to those in the population (Gregory et al., 2005;
117 Van Swaay et al., 2002, 2008; Weiser et al., 2020). But it is also common to see samples
118 treated as though they are representative despite clear evidence to the contrary. For example,
119 Vellend et al. (2013) and Dornelas et al. (2014) purported to document globally
120 representative time trends in species richness, but Gonzalez et al. (2016) showed that their
121 samples were highly unrepresentative with respect drivers of biodiversity change and species
122 richness itself. (See Boyd, Powney, et al. (2023) for a review of this debate and others like it.)
123 We suspect that many of those who do not deal with issues of sample representativeness are
124 not familiar with the gravity of the problem or the relevant theory and adjustment methods.

125 In this paper, we introduce six approaches to descriptive inference using unrepresentative
126 nonprobability samples and demonstrate how they relate to each other (conceptually and
127 mathematically). We apply each approach to a simple two-part biodiversity monitoring
128 problem. The first part is to estimate mean occupancy of the plant *C. vulgaris* across 1 km
129 grid squares in Britain in two time-periods; the second is to estimate the difference between
130 the two (i.e. the time trend). *Calluna vulgaris* is an attractive case study because we have
131 good estimates of its true geographic distribution in both periods from satellite (amongst
132 other sources). The approaches to inference that we demonstrate are subsampling, quasi-
133 randomisation (Elliott and Valliant, 2017), poststratification (Little, 1993), superpopulation
134 modelling (Valliant, 2009), a “doubly robust” estimator (Chen et al., 2020), and Multilevel
135 Regression and Poststratification (MRP; Gelman, 2007; Gelman and Little, 1997). Each can
136 be (MRP more loosely than the rest) interpreted as an attempt to weight the sample in such a
137 way that it more closely resembles the population, in the hope that this results in more
138 accurate descriptive inferences. We demonstrate the effects of each approach on the
139 distributions of auxiliary variables in the sample, as well as on the resulting estimates of
140 mean occupancy in each period and the time trend between the two. Applying the adjustment

141 methods to a real-world example reveals challenges that ecologists are likely to face, and we
142 discuss these in detail.

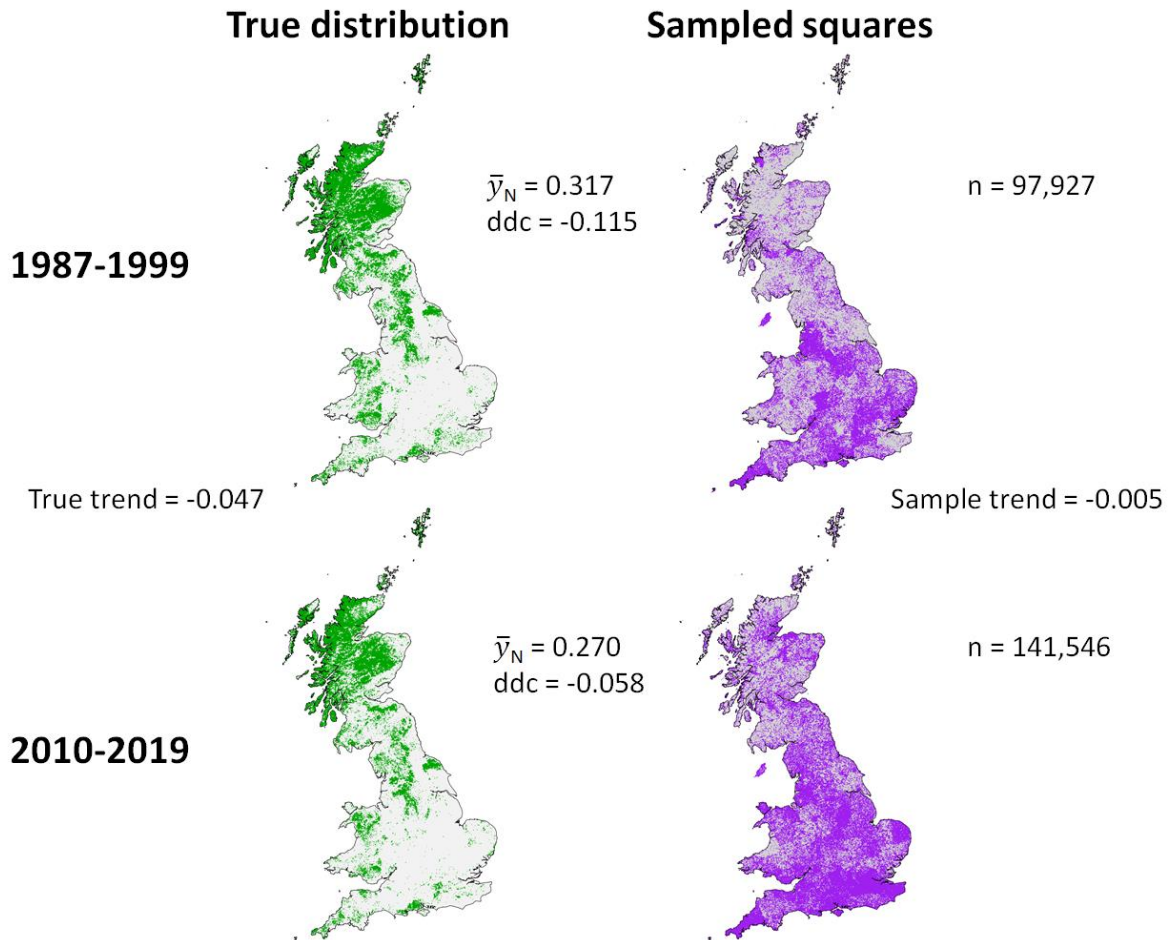
143 **Methods**

144 **True distribution of *Calluna vulgaris***

145 We approximated the true distribution of the dwarf shrub vascular plant *Calluna vulgaris*
146 (Heather) in two time periods: 1987–1999 and 2010–2019. For the first period, we used the
147 1990 UKCEH land cover map (Rowland et al., 2020); for the second, we used the 2018
148 version (Morton et al., 2022). The land cover maps are derived from satellite, which means
149 that they provide information for every 1 km grid square. From these maps, we identified 1
150 km grid squares (British National Grid, EPSG:27700) with >0% heather or heather grassland
151 cover. To these, we added 1 km squares in which *C. vulgaris* was recorded in each time
152 period by the Botanical Society of Britain and Ireland (BSBI). The time periods considered
153 cover the main periods of recording for two national distribution atlases, which involved a
154 concerted effort by volunteers (citizen scientists) to document vascular plants across the
155 United Kingdom (Preston, C.D., Pearman, D.A. & Dines, 2002; Stroh et al., 2023).

156 Acknowledging that some 1 km squares may have been erroneously classed as having some
157 heather or heather grassland coverage by the land cover maps, we removed any 1 km squares
158 that fell within 10 km grid squares in which *C. vulgaris* had not been recorded by the BSBI in
159 the period 1950–2019. Given that this period includes recording for three national
160 distribution atlases (the two cited above plus Perring & Walters, 1962), we assume that the
161 union of all 10 km occurrences within this period encompasses all known populations
162 irrespective of finer scale changes. Figure 1 maps the resulting estimates of the true 1 km
163 distributions of *C. vulgaris* in both time-periods.

164



165

166 Figure 1. Left column: the distribution of *Calluna vulgaris* in both time-periods. Green
 167 squares are occupied and grey squares are not. \bar{y}_N is mean occupancy or, equivalently, the
 168 proportion of squares occupied. The ddc's are the correlations between sample inclusion (1 if
 169 the square is in the sample and 0 otherwise) and occupancy. Right column: the nonprobability
 170 1 km samples for each time-period. Purple squares were sampled and grey squares were not.
 171 n is the number of squares sampled. We assume that *C. vulgaris* was recorded in all sampled
 172 grid squares that it occupied in the relevant time-period. The true trend is the difference
 173 between population means, and the sample trend is the difference between sample means (i.e.
 174 mean occupancy across purple squares).

175 **Sample data on *Calluna vulgaris* occupancy**

176 The 1 km samples for both time periods ("Sampled squares in Fig. 1) encompass any vascular
 177 plant data for which the date of collection is known (i.e. the record is resolved to the day),
 178 either at the 1 km scale or finer, collected by the BSBI for the national distribution atlases of
 179 Preston et al. (2002) and Stroh et al. (2023). Having been collected by volunteers, the data
 180 come under the banner of citizen science.

181 **Auxiliary data**

182 We used two auxiliary variables for which data are available for all 1 km grid squares in
 183 Great Britain: the proportion of each 1 km grid square that falls within some form of
 184 protected area (including everything from SSSI's to local nature reserves; UNEP-WCMC &
 185 IUCN, 2020) and the average elevation of each 1 km square (Intermap, 2009). New protected
 186 areas are designated periodically, so we used the set that were designated prior to 1987 for
 187 the first time-period and prior to 2010 for the second (i.e. the beginning of each period). We

188 suspect that 1 km squares with more protected area coverage are more likely to be visited by
189 naturalists (Girardello et al., 2019) and, because protected areas tend to have higher quality
190 habitat, are also more likely to be occupied by *C. vulgaris*. Likewise, elevation should affect
191 both sample inclusion and *C. vulgaris* occupancy. Sites at higher elevations are harder to
192 access on account of their relatively harsh terrain and remoteness, and elevation is a known
193 predictor of *C. vulgaris* occupancy (Stroh et al., 2023).

194 One of the adjustment methods that we describe below, quasi-randomisation, requires
195 additional covariates (we use the term “covariate” to distinguish these from the auxiliary
196 variables as defined earlier). The method involves the estimation of sample inclusion
197 probabilities for every 1 km grid square in Britain. This is a matter of prediction rather than
198 inference, because we know whether each 1 km square was sampled (i.e. there is no missing
199 data), so it was sensible to use a wider range of covariates. See Table 1 in Boyd, Stewart, et
200 al. (2023) for a list of the additional covariates used in this model.

201 **Estimating the per-period population mean**

202 The first step in our biodiversity monitoring problem is to estimate mean occupancy of *C.*
203 *vulgaris* in each time-period. Although not usually written this way, it is helpful for what
204 comes later to re-express the population mean as a weighted sum

$$\bar{y}_N = \frac{1}{N} \sum_{i=1}^N y_i = \sum_{i=1}^N \frac{y_i}{N} = \sum_{i=1}^N \frac{y_i w_i}{\sum_{i=1}^N w_i}, \quad \text{equation 1}$$

205

206 where y is occupancy (1 = occupied and 0 = unoccupied), N is the population size, i indexes
207 1 km grid squares and $w_i = 1/N$ (N is the same in both time-periods). The denominator in
208 the rightmost expression might seem unnecessary, because it equals one. We have retained it
209 to illustrate the similarity between this expression and the sample-based estimators below,
210 which have a similar form but whose sampling weights w do not necessarily sum to one. (We
211 use the term “estimator” to describe a rule for estimating some quantity from a sample; here,
212 that quantity is the population mean.) For notational simplicity, we do not index the time-
213 period, and the reader should remember that \bar{y}_N is time-period specific. In practice, y is not
214 known for all i in the population, so sample-based estimators of \bar{y}_N are needed.

215 **The design-based estimator**

216 The design-based estimator of the population mean, which is applicable only where a
217 probability sample of some sort is available (Lohr, 2022), has a similar form to eq. 1

$$\bar{y}_{ab} = \sum_{i=1}^n \frac{y_i w_i}{\sum_{i=1}^n w_i}. \quad \text{equation 2}$$

218 The differences are that the sums are over the sample size n rather than N and that the
219 weights w_i are not necessarily constant. Rather, the weight for unit i , w_i , is equal to the
220 reciprocal of the probability that it was included in the sample = $1/p_i$.

221 Sample inclusion probabilities are, by definition, not known for nonprobability samples, so
222 alternative estimators are required. We present six such estimators below, three of which—
223 quasi-randomisation, poststratification and superpopulation modelling—are explicit attempts
224 to come up with a set of weights w_i that produce a reasonable estimate of \bar{y}_N under eq. 2. The
225 other three—a “doubly robust” estimator, subsampling and MRP—are not, but they are
226 conceptually similar.

227 **Estimators for nonprobability samples**

228 The following estimators are used in survey sampling to estimate population means from
229 nonprobability samples. More detail on each can be found in Valliant et al. (2018), Lumley
230 (2010) and Lohr (2022). See supplementary material 1 for an R Markdown document
231 containing the code to implement each of the adjustment methods.

232 ***Naïve sample mean***

233 Where sample inclusion probabilities are unavailable, a simple option is to assume that $w_i =$
234 $1/n$ for all i . In this case, eq. 2 gives the (naïve) sample mean. As the weights are constant,
235 the sample mean does not adjust for differences in y between the sampled and non-sampled
236 population units. It is nevertheless widely used in biodiversity monitoring.

237 ***Quasi-randomisation***

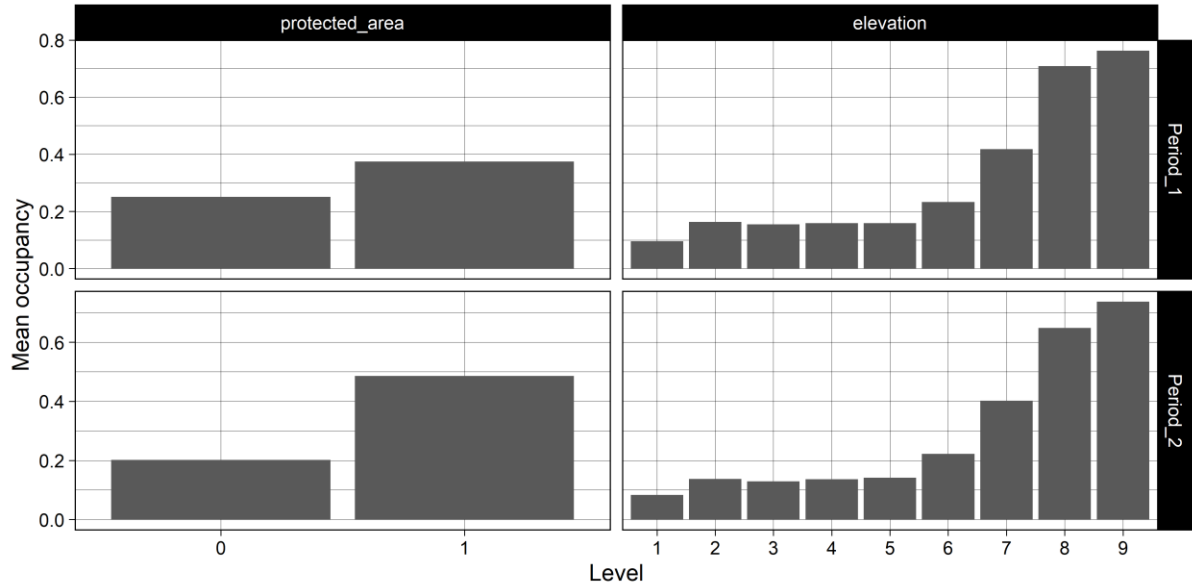
238 An alternative approach is to imagine that the nonprobability sample was selected
239 probabilistically and to estimate the implied inclusion probabilities. Any binary model and
240 covariates can be used. Once inclusion probabilities p_i have been estimated, the weights $w_i =$
241 $1/p_i$ (as in the design-based estimator). In our example, we used random forests and several
242 covariates (including the auxiliaries) to estimate pseudo-inclusion probabilities. More
243 complex approaches are possible and have been used to map species distributions (Johnston
244 et al., 2020).

245 ***Poststratification***

246 Another approach to estimating sampling weights is poststratification. Poststratification
247 requires categorical auxiliary data, so continuous variables must be discretized prior to
248 analysis (Valliant, 2020). The auxiliary variables are crossed (think contingency tables) to
249 create poststrata. Each poststratum j has a sample size n_j and population size N_j . The
250 sampling weight w_i for population unit i in poststratum j is given by N_j/n_j .

251 In our example, we split elevation into ten categories using its deciles (i.e. cut points at the
252 10th and 20th percentiles, etc.). This did not make sense for the variables denoting the
253 proportion of each grid square that falls within a protected area, because most squares took
254 the value one or zero. We split this variables into two categories, 0 and >0 , i.e. whether or not
255 there is some protected land in the grid square. Discretization gave $10 \times 2 = 20$ poststrata.

256 It is sensible to discretize the auxiliary variables in such a way that the variable of interest
257 varies among categories. Otherwise, the adjustment from poststratifying will be minor (or
258 unnecessary!). Fig. 2 shows that mean occupancy of *C. vulgaris* in the samples differs
259 appreciably among levels of the auxiliary variables.



260

261 Figure 2. Mean occupancy of *Calluna vulgaris* for each level of the auxiliary variables in
 262 each time-period. The auxiliary variables were originally on a continuous scale, but we
 263 discretized them to enable poststratification. See the main text for details.

264 **Superpopulation modelling**

265 Superpopulation modelling is conceptually different to the adjustment methods described
 266 above. The premise is that there exists some model that describes the variable of interest in
 267 the population. If this model can be recovered from the sampled outcome variable y and the
 268 auxiliary data, it can be used to predict the variable of interest in non-sampled units. Given a
 269 prediction for each non-sampled i , it is then simple to estimate the population mean.

270 A general (i.e. multiple) linear regression model of y has the form

$$E_M(y_i) = \mathbf{x}_i^T \boldsymbol{\beta}, \quad \text{equation 3}$$

271 where the subscript M indicates that the expectation (mean) is with respect to the model, \mathbf{x}_i is
 272 the vector of auxiliary variables for unit i , the superscript T indicates that the vector \mathbf{x}_i has
 273 been transposed (to a row vector) and $\boldsymbol{\beta}$ is a column vector of parameters. A prediction of y
 274 for unit i is

$$\hat{y}_i = \mathbf{x}_i^T \hat{\boldsymbol{\beta}}. \quad \text{equation 4}$$

275 The accent on $\hat{\boldsymbol{\beta}}$ indicates that it is an estimate (the least squares estimate in this case). If \bar{s} is
 276 the set of non-sampled population units, the superpopulation model prediction of the
 277 population mean is

$$\bar{y}_{sp} = \frac{\sum_{i \in s} y + \sum_{i \in \bar{s}} \hat{y}}{N}. \quad \text{equation 5}$$

278 That is, it is the sum of the known outcome values in the sample and those predicted by the
 279 model for the remainder of the population divided by the population total.

280 A feature of \bar{y}_{sp} is that it can be expressed in the same form as the design-based estimator in
 281 eq. 2, with the weights w_i being a function of the auxiliary variables in sampled and non-
 282 sampled population units (Elliott & Valliant, 2017). (Code to verify this numerically is
 283 available at <https://github.com/robboyd/selectionBiasEffects/tree/master/R>.) Like the other

284 adjustment models, then, the superpopulation estimator is an approach to estimating the
285 sampling weights w_i .

286 Linear regression might seem like an unusual choice of model for a binary outcome
287 (occupancy), but we felt that it was the best option here. One reason is that the implied model
288 is actually linear for an estimator with the form of eq. 2 (Valliant, 2020). Most important,
289 however, is that the use of a linear model enables the estimation of sampling weights
290 (Valliant et al., 2018; supplementary material 1;
291 <https://github.com/robboyd/biasAdjustments>). This is helpful, because those weights can be
292 used to show the effects of superpopulation modelling on the distributions of the auxiliary
293 variables in the sample (see “Evaluating the effects of the adjustments” below). Alternative
294 models can be used where weights are not required (e.g. Wu and Sitter, 2001). In our
295 example superpopulation model, we used the auxiliary variables as predictors.

296 ***Doubly robust estimator***

297 The doubly robust estimator combines the superpopulation model and the sample inclusion
298 model from the quasi-randomisation procedure in such a way that if either is correct, and the
299 sample size is large, then the estimate of the population mean unbiased (Valliant, 2020). It
300 has the general form (Wu, 2022)

$$\bar{y}_{dr} = \frac{1}{N} \sum_{i \in s} \frac{r_i}{p_i} + \frac{1}{N} \sum_{i=1}^N \hat{y}_i, \quad \text{equation 6}$$

301 where $r_i = y_i - \hat{y}_i$ (i.e. the residuals of superpopulation model). The second term on the right
302 is the superpopulation model prediction of \bar{Y}_N . If the superpopulation model is correctly
303 specified, then it is an unbiased estimate of \bar{Y}_N . However, if the superpopulation model is
304 misspecified, then the second term needs to be corrected, which is where the first term comes
305 in. If the quasi-randomisation sample inclusion model is correctly specified, the first term
306 corrects the second by adding the residuals of the superpopulation model divided by the
307 (correctly) estimated pseudo inclusion probabilities. This is sufficient to produce an unbiased
308 estimate of \bar{Y}_N even where the superpopulation model is wrong. Where the superpopulation
309 model is correct, the first term is 0, because $r_i = 0$. Where neither model is correct, \bar{y}_{dr} is a
310 biased estimator of \bar{Y}_N . See Chen et al. (2020), who combined probability and nonprobability
311 samples, for a similar approach.

312 ***Subsampling***

313 Perhaps more familiar to ecologists than the above approaches is subsampling (Beck et al.,
314 2014; Steen et al., 2020). The idea is to create a representative “miniature” of the population
315 out of the sample (Meng, 2022) and to calculate the quantity of interest (mean occupancy)
316 from this subsample. Subsampling trades sample size for representativeness.

317 Our approach was to draw stratified random samples of size $N/10 = 22,958$ with
318 replacement from the original samples. We used the same strata as described above (under
319 Poststratification). The decision to set $n = N/10$ was somewhat arbitrary, but changing the
320 subsample size makes little difference to the point estimates of the population means
321 (although they become more precise with increasing subsample size; supplementary material
322 1). The subsample mean is the estimator of the population mean.

323 ***Multilevel regression and poststratification (MRP)***

324 MRP is an extension of poststratification and a variation of superpopulation modelling
325 (Gelman, 2007; Gelman & Little, 1997; Valliant et al., 2018). A hierarchical model is used to
326 estimate mean occupancy in each poststratum. The advantage of using a hierarchical model is
327 that cells with few or no data borrow information from cells with more data (i.e. partial

328 pooling or shrinkage is exploited). The population mean is the weighted mean of the stratum
329 means, where the weights are equal to the proportion of the population in each stratum.

330 Our hierarchical model is a binomial GLM with a logit link function, a fixed intercept and
331 random intercepts for the auxiliary variables and their interaction (see [https://mc-](https://mc-stan.org/rstanarm/articles/mrp.html)
332 [stan.org/rstanarm/articles/mrp.html](https://mc-stan.org/rstanarm/articles/mrp.html) for a similar formulation). We fitted the model in a
333 Bayesian framework using 5 Markov Chain Monte Carlo (MCMC) chains, each with 1000
334 iterations. This was sufficient to achieve convergence on all parameters in both time-periods.

335 **Confidence intervals**

336 We present 95% confidence/credible intervals for all estimates of mean occupancy (credible
337 intervals for MRP, which we implemented in a Bayesian framework). For most methods—
338 superpopulation modelling, quasi-randomisation, subsampling and the doubly robust
339 estimator—we constructed bootstrap confidence intervals. Resampling the original data with
340 replacement, we created 1000 bootstrap samples, from which we obtained a distribution of
341 estimates from each method and calculated percentile intervals. For MRP, we extracted
342 credible intervals from the posterior distributions of mean occupancy. We used the
343 confidence intervals provided by the *survey* package (Lumley, 2010) for the poststratified and
344 naïve (i.e. unadjusted) estimates.

345 **Estimating the trend in mean occupancy**

346 Having estimated mean occupancy in each time-period, the next step was to estimate the
347 difference between the two = $\bar{y}_2 - \bar{y}_1$ (i.e. the trend). We constructed a confidence interval
348 for the trend estimated using each method in one of two ways depending on whether the
349 method produced one estimate or a distribution. The methods that produced a distribution of
350 $\bar{y}_2 - \bar{y}_1$ include those that we bootstrapped and MRP, which we fitted in a Bayesian
351 framework (meaning we have a posterior distribution). For these methods, we extracted
352 percentile confidence intervals (95%) from the distributions of estimated trends. For the
353 others, poststratification and the naïve estimator (the sample mean), we used the normal
354 approximation of the 95% confidence interval, given by $\pm 1.96 \times$ the standard errors, where
355 the standard errors are $\sqrt{\text{var}(\bar{y}_2) + \text{var}(\bar{y}_1)}$ (Gelman, 2007).

356 **Evaluating the effects of the adjustments**

357 We used relative frequency plots (Cf. Makela et al., 2014) to assess whether the adjustments
358 brought the distributions of the auxiliary variables in the samples closer to their distributions
359 in the population. The first step was to split each auxiliary variable into fifty bins of equal
360 width spanning its range. The relative frequency of grid squares (the i 's) in each bin k is
361 $N_{i,k}/N$, where $N_{i,k}$ is the number of grid squares in each bin k in the population and N is the
362 population size (we use k to index the bins to distinguish them from the strata described
363 earlier). Similarly, the relative frequency of sampled grid squares in each k is $n_{i,k}/n$, where
364 $n_{i,k}$ is the number of sampled grid squares in bin k and n is the total sample size. In the
365 adjusted samples, the equivalent relative frequency is $\frac{\sum_{i \in k} w_i}{\sum_N w_i}$ (slightly different for
366 subsampling; see below). We compared the original and adjusted samples' deviations from
367 the population using the Mean Absolute Error (MAE) of the relative frequencies across all k .
368 If the MAE from the adjusted sample is smaller than the original sample, then the adjustment
369 brought the distribution of the auxiliary variable closer to its population distribution.

370 We were not able to produce adjusted relative frequency plots based on the doubly robust
371 estimator or MRP. The problem was that could not estimate reasonable sampling weights
372 from either method, which are needed to adjust the relative frequencies of the auxiliaries.
373 Whilst it has been shown how to derive unit-level sampling weights where the MRP

374 multilevel model is linear (Gelman, 2007), no formula has yet been derived for the case of the
 375 binomial GLM (Valliant et al., 2018). As for the doubly robust estimator, Valliant (2020)
 376 showed how to derive “model-assisted” weights. Unfortunately, in our case, many of the
 377 model-assisted weights were very large and negative. The extreme weights appear to be
 378 caused by the pattern of residuals from the superpopulation model (recalling that we used a
 379 linear regression despite the fact that occupancy is binary), but it is beyond the scope of this
 380 paper to definitively diagnose the problem. There is no obvious way to derive weights from
 381 the subsampling estimator either. However, for this estimator, the adjusted relative
 382 frequencies of the auxiliaries are simply their distributions in the subsamples so are simple to
 383 obtain.

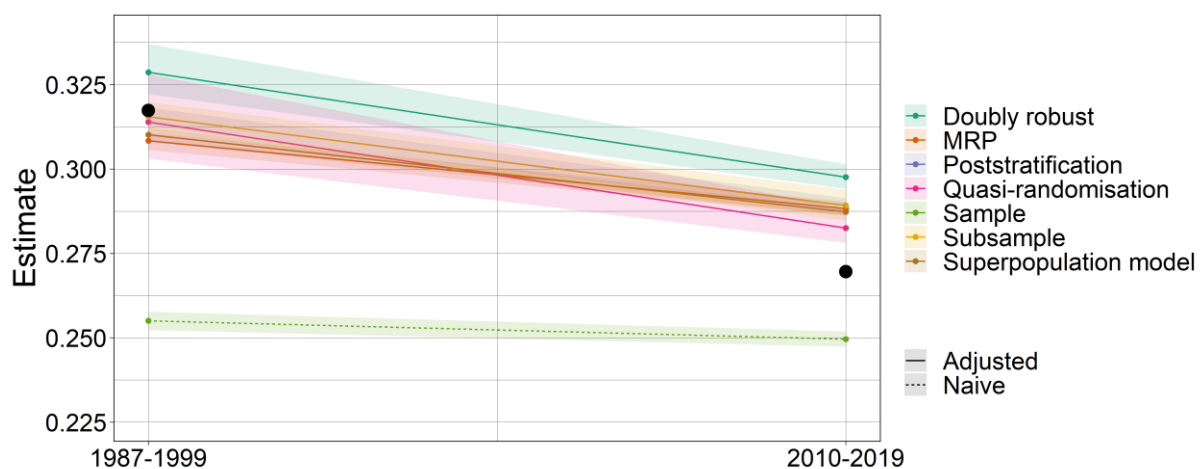
384 Assessing whether the estimates of mean occupancy in each period and the trend were
 385 improved by each adjustment method was simpler. We measured the difference between the
 386 point estimates of mean occupancy and the truth using the absolute error = $|\bar{y}_N - \bar{y}_{est}|$, where
 387 \bar{y}_{est} is the estimate. For the trends, whose signs are of interest, we simply used the
 388 differences between the estimates and the truth. We also assessed whether the
 389 confidence/credible intervals produced by each method covered the true means and trend. We
 390 did not consider the power to detect the trend—that is, whether the methods’ uncertainty
 391 intervals span zero at some percentile—because many biodiversity applications are
 392 descriptive-inferential rather than decision-theoretic.

393 Results

394 Per-period sample representativeness and estimated mean occupancy

395 The samples are large but somewhat unrepresentative (Fig. 1). Forty-three percent of grid
 396 squares were sampled in period one, and the ddc is -0.115; in period two, 62% of grid squares
 397 were sampled, and the ddc is -0.057. A consequence of these ddc’s is that the naive sample
 398 means underestimate the population means, especially in period one where the magnitude of
 399 the ddc is greater (Fig. 3).

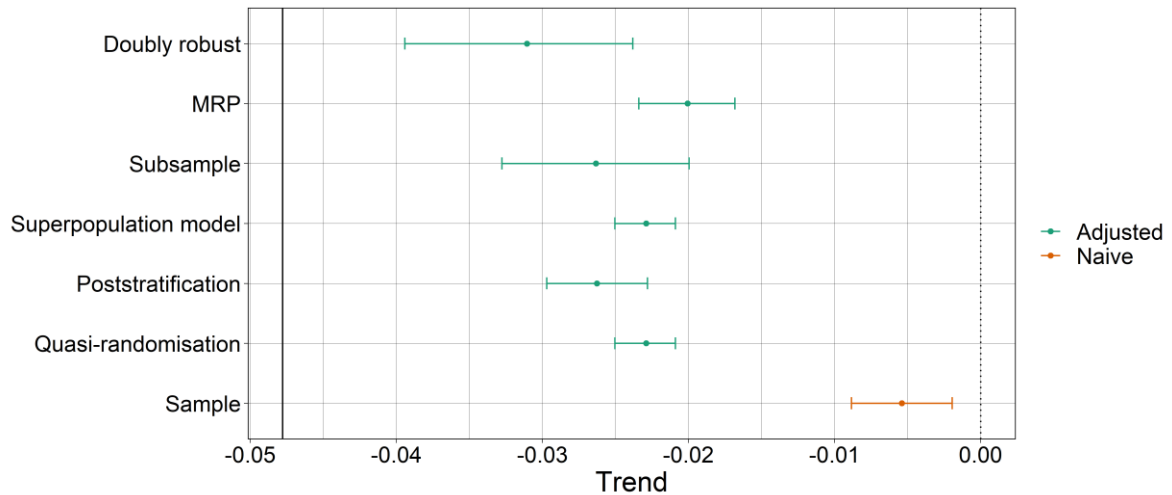
400 With the exception of the doubly robust estimate in period two, the estimates of mean
 401 occupancy from all adjustment methods in both time-periods had lower absolute errors than
 402 the naïve sample mean (Fig. 3; mean absolute errors are provided in supplementary material
 403 2). The confidence intervals for the poststratified, subsample and quasi-randomisation
 404 estimates covered the true population mean in period one. In period two, no method’s
 405 confidence/credible interval covered the population mean.



407 Figure 3. Naive (i.e. unadjusted) and adjusted sample-based estimates of mean occupancy in
 408 each time-period. The shaded regions are 95% confidence/credible intervals (see the main
 409 text for information on we constructed these for each method). The large black circles are the
 410 true population means in each time-period.

411 **Estimated trend in mean occupancy**

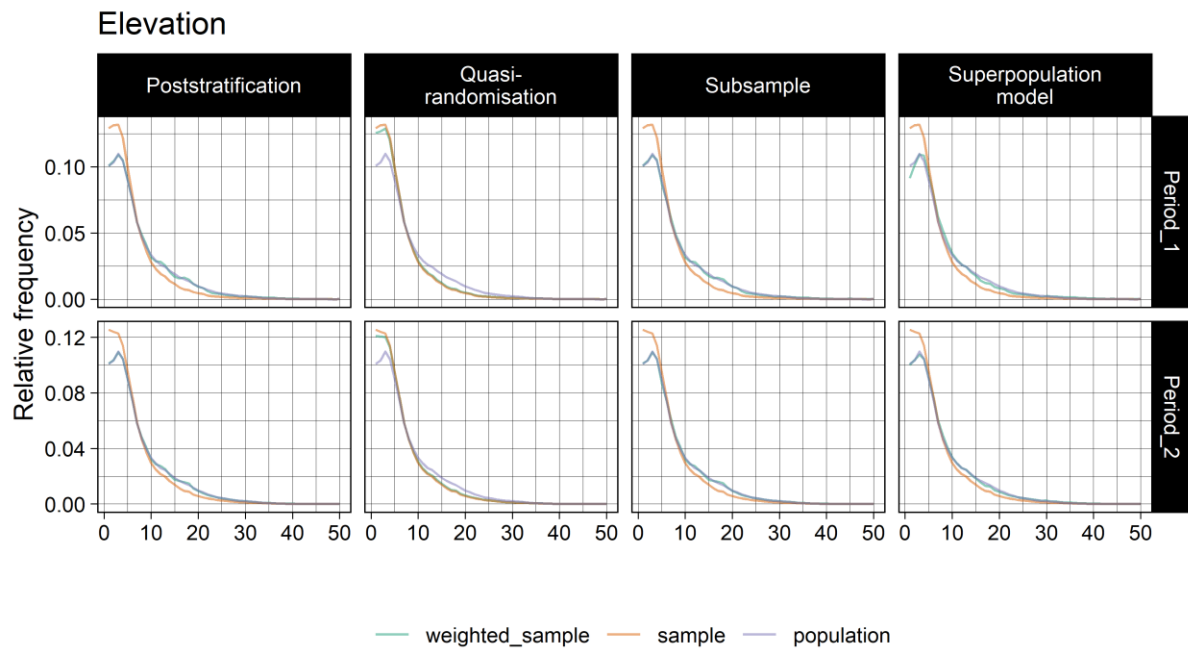
412 Estimates of the trend in mean occupancy from all adjustment methods were more accurate
 413 than the difference in sample means (i.e. the naive estimate; Fig. 4). However, no method's
 414 point estimate came close to the true trend of -0.047, and their confidence/credible intervals
 415 did not cover it.



416
 417 Figure 4. Trends in mean occupancy between periods one and two produced by the estimator
 418 from each adjustment method, in addition to the naive sample estimate. Error bars delimit
 419 95% confidence/credible intervals. The solid vertical black line denotes the true population
 420 trend (-0.047).

421 **Distributions of auxiliary variables**

422 As measured using Mean Absolute Errors (MAEs), the adjustment methods were generally
 423 very good at bringing the distributions of the auxiliaries in the samples closer to those in the
 424 population. Fig 5 shows the sample and population distributions of elevation, but the MAEs
 425 for this and the proportion of each grid square that falls within a protected area can be found
 426 in supplementary material 2.



427 Figure 5. Sample, population and weighted sample distributions of the auxiliary variable road
 428 length (Table 1) in periods one and two.
 429

430 **Discussion**

431 We applied six approaches to descriptive inference from nonprobability samples to a simple
 432 biodiversity monitoring problem: the estimation of mean occupancy of the plant *C. vulgaris*
 433 in two time-periods and the trend between the two. The methods generally worked well in the
 434 sense that they brought the distributions of auxiliary variables in the samples closer to their
 435 distributions in the population (all 1 km grid squares in Britain). Successful redistribution of
 436 the auxiliaries translated into improvements of the estimates of mean occupancy in both time-
 437 periods and the trend between the two. Importantly, however, no method was completely
 438 unbiased, and their uncertainty intervals did not cover the true values of occupancy in the
 439 second period or the trend. An abatement rather than an elimination of bias is probably the
 440 best outcome that can be expected, because most adjustment methods rest on the untenable
 441 assumption that non-sampled locations are “Missing At Random” (MAR); that is, the
 442 variable of interest is completely independent of sample inclusion given the auxiliary
 443 variables.

444 Unlike most practical situations, we were able to test the MAR assumption, because we know
 445 the true distribution of *C. vulgaris* in Britain. In the first time-period, the partial correlation
 446 between sample inclusion and occupancy, conditional on elevation and protected area
 447 coverage, is -0.018; in period two, it is 0.035 (supplementary material 1). These “adjusted”
 448 ddc’s are lower in magnitude than the original ddc’s, -0.115 and -0.058, which means that
 449 accounting for elevation and protected area coverage increased the representativeness of the
 450 samples (recalling that a smaller ddc means a more representative sample). That is not to say
 451 that the samples became *fully* representative, which would be the case in expectation in a
 452 MAR scenario. The usual yardstick for a representative sample is the simple random sample,
 453 whose ddc is of the order $N^{-1/2}$ (Meng, 2018). In our example, $N^{-1/2} = 2.2^{-6}$, which is
 454 several orders of magnitude smaller than the “adjusted” ddc’s. This goes to show that without
 455 a truly miniscule ddc, which would only be induced (in expectation) where the MAR
 456 assumption holds or under random sampling, sample means as estimators of population
 457 means will be appreciably biased (especially where N is large).

458 It might seem wise to include as many potential auxiliaries as possible to reduce the chance
459 of missing a genuine one. For example, Collins et al. (2001) advocated for including all
460 variables exceeding some prescribed correlation with sample inclusion and the variable of
461 interest. This strategy can be a dangerous one, however. Thoemmes & Rose (2014) show that
462 including correlates of sample inclusion and the variable of interest, rather than theoretically
463 justifiable causes, can increase the bias in estimates of population means (also see Thoemmes
464 & Mohan, 2015). Indeed, in a previous version of this manuscript (Boyd, Stewart, et al.,
465 2023), we took a more inclusive approach to the selection of auxiliary variables, and our
466 estimates of *C. vulgaris* occupancy in period two were generally more biased than the naïve
467 estimate from the unadjusted sample.

468 Identifying appropriate auxiliary variables is likely to be the most challenging part of
469 adjusting samples in biodiversity monitoring. In many situations, causes of the variable of
470 interest and sample inclusion are not known. Taxon and dataset experts might be able to
471 identify potential auxiliary variables, but it is unlikely that they can identify them all (which
472 would be needed to satisfy the MAR assumption). The experts might also erroneously
473 identify auxiliary variables that are not suitable, in which case adjusting for those variables
474 might do more harm than good (Thoemmes & Rose, 2014). Even if experts were able to
475 correctly identify all relevant auxiliaries, those variables might not be reflected in available
476 data. Transparency regarding availability and choice of auxiliary variables should be an
477 important component of reporting for all biodiversity monitoring.

478 Acknowledging that variables of interest in biodiversity monitoring are likely to be dependent
479 on sample inclusion even after controlling for the available auxiliaries, it might be worth
480 considering adjustment methods that forgo the MAR assumption. For example, Tchetgen
481 Tchetgen & Wirth (2017) showed it is possible to recover a true population regression model
482 (and therefore the population mean) by incorporating “instrumental variables”. They define
483 instrumental variables as those that are predictive of sample inclusion, independent of the
484 variable of interest and independent of “selection bias” (the latter defined as the mean of the
485 variable of interest in the sample minus the mean of the variable of interest in non-sampled
486 population units). We screened three additional variables—the proportion of each grid square
487 that is accessible to the public, the density of postcodes in each grid square and its nearest
488 neighbours, and the length of major roads in each grid square and its nearest neighbours—to
489 see if they satisfied these three assumptions, but none did (supplementary material 1). In
490 practical situations, where the variable of interest is not known for non-sampled population
491 units, testing these assumptions would be challenging.

492 Whilst we are confident that the availability of data on auxiliary variables was the limiting
493 factor in our example, it is possible that improvements to the adjustment methods themselves
494 could have improved matters. Where sampling weights are not of interest, for example, it
495 might be sensible to use a binomial generalised linear model, rather than a general linear
496 regression, for the superpopulation model (Wu & Sitter, 2001). The multilevel modelling
497 component of MRP exploits partial pooling, so we could have used more finely resolved
498 strata on the basis that estimates for sparse strata (with low sample sizes) would be shrunk
499 towards those from strata with more data. The question is whether fine-tuning the adjustment
500 methods is likely to result in large improvements in accuracy. As Mercer et al., (2018),
501 writing in the context of adjusting survey samples, put it, “[t]he right variables make a big
502 difference for accuracy. Complex statistical methods, not so much.” The fact that most
503 adjustment methods performed almost similarly in our example is further evidence that the
504 choice of auxiliary variables matters more than the specifics of the adjustment method.

505 Given that the methods performed similarly in terms of accuracy, it would be sensible to
506 consider those that are quickest to run. As we implemented it, MRP took by far the longest to
507 run of all the methods—about ten hours per time-period on a computer cluster. Bootstrapping
508 to estimate confidence intervals meant that other methods, too, were quite expensive to run.
509 This was particularly true for the quasi-randomisation and doubly robust procedures, both of
510 which involved repeatedly fitting the sample inclusion model—itself a time consuming
511 process. The remainder of the methods—superpopulation modelling, subsampling and
512 poststratification—took a negligible amount of time to run.

513 Although we have only considered one species and dataset, previous studies (in other
514 disciplines) shed light on the factors that affect the accuracy of inference from nonprobability
515 samples more generally. Omitting genuine auxiliary variables in the adjustment process is
516 more problematic where those variables explain larger proportions of the variance in the
517 variable of interest and sample inclusion (Collins et al., 2001). Equally, inclusion of certain
518 variables that are not appropriate auxiliaries becomes more problematic where they explain
519 larger proportions of the variance in the variable of interest and sample inclusion (Thoemmes
520 & Rose, 2014). In practice, we do not know the strengths of the effects potential auxiliaries
521 on the variable of interest and sample inclusion, or whether they have effects at all, but it is
522 clear that the selection of auxiliary variables will be a critical component of adjusting samples
523 in biodiversity monitoring.

524 Given the importance of selecting appropriate auxiliary variables, we propose the following
525 general strategy for analysts intending to draw inferences about biodiversity change from
526 geographically unrepresentative nonprobability samples. The first step should be to consult
527 taxon and dataset experts, who might be able to identify relevant auxiliary variables. Where
528 possible, consulting multiple experts to capture their uncertainty about what affects sample
529 inclusion and the variable of interest would be desirable. If data are available on these
530 variables, then their distributions in the sample and population should be compared to assess
531 whether the data are representative with respect to that variable. Several tools are available to
532 perform such comparisons (Boyd et al., 2021; Ruete, 2015). The next step should be to adjust
533 the sample based on the relevant auxiliaries and to draw inferences from the adjusted
534 samples. Like others (e.g. Mercer et al., 2018), we found that it is of little consequence which
535 adjustment method is used, so it is sensible to pick one that is quick to run. Rather than
536 assuming the adjustment worked perfectly, it is important to acknowledge and report the
537 potential for residual bias. As we have shown, traditional uncertainty intervals are not
538 guaranteed (or even likely) to cover the true population parameters of interest unless all
539 relevant auxiliaries are known and reflected in available data (Meng, 2018). Where there is
540 doubt about the relevant auxiliary variables, a safer strategy is to assess the risk of bias
541 qualitatively and to ensure it is reflected in the way that findings are reported (Boyd et al.,
542 2022; Meineke & Daru, 2021; Pescott et al., 2022).

543

544 **Acknowledgements**

545 Thank you to Richard Valliant and two anonymous reviewers, whose comments improved
546 this paper. All authors were supported by the NERC Exploring the Frontiers award number
547 NE/X010384/1 “Biodiversity indicators from nonprobability samples: Interdisciplinary
548 learning for science and society”. OLP was also supported by the NERC award number
549 NE/R016429/1 as part of the UK Status, Change and Projections of the Environment (UK-
550 SCAPE) programme delivering National Capability.

551 **References**

- 552 Bailey, M. (2022). Comments on “Statistical inference with non-probability survey samples
553 .” *Survey Methodology*, 48(12), 331–338.
- 554 Bethlehem, J., Cobben, F., & Schouten, B. (2008). Indicators for the Representativeness of
555 Survey Response. *Statistics Canada’s International Symposium Series: Proceedings*, 11.
- 556 Bowler, D. E., Bhandari, N., Repke, L., Beuthner, C., Callaghan, C. T., Eichenberg, D.,
557 Henle, K., Klenke, R., Richter, A., Jansen, F., Bruelheide, H., & Bonn, A. (2022).
558 Decision-making of citizen scientists when recording species observations. *Scientific*
559 *Reports*, 12(1), 1–12. <https://doi.org/10.1038/s41598-022-15218-2>
- 560 Bowler, D. E., Klaus, D. E., Conze, J., Suhling, F., Baumann, K., Benken, T., Bönsel, A.,
561 Bittner, T., Drews, A., Günther, A., Isaac, N. J. B., & Petzold, F. (2021). Winners and
562 losers over 35 years of dragonfly and damselfly distributional change in Germany.
563 *Diversity and Distributions*, 27(August 2020), 1353–1366.
564 <https://doi.org/10.1111/ddi.13274>
- 565 Boyd, R. J., Powney, G., Carvell, C., & Pescott, O. L. (2021). occAssess: An R package for
566 assessing potential biases in species occurrence data. *Ecology and Evolution*,
567 11(September), 16177–16187. <https://doi.org/10.1002/ece3.8299>
- 568 Boyd, R. J., Powney, G. D., Burns, F., Danet, A., Duchenne, F., Grainger, M. J., Jarvis, S. G.,
569 Martin, G., Nilsen, E. B., Porcher, E., Stewart, G. B., Wilson, O. J., & Pescott, O. L.
570 (2022). ROBITT: A tool for assessing the risk-of-bias in studies of temporal trends in
571 ecology. *Methods in Ecology and Evolution*, 13(March), 1497–1507.
572 <https://doi.org/10.1111/2041-210X.13857>
- 573 Boyd, R. J., Powney, G. D., & Pescott, O. L. (2023). We need to talk about nonprobability
574 samples. *Trends in Ecology & Evolution*, xx(xx), 1–11.
575 <https://doi.org/10.1016/j.tree.2023.01.001>
- 576 Boyd, R. J., Stewart, G. B., & Pescott, O. L. (2023). Descriptive inference using large,
577 unrepresentative nonprobability samples: An introduction for ecologists. *Ecoevorxiv*,
578 April. <https://doi.org/10.32942/X2359P>
- 579 Chen, Y., Li, P., & Wu, C. (2020). Doubly Robust Inference With Nonprobability Survey
580 Samples. *Journal of the American Statistical Association*, 115(532), 2011–2021.
581 <https://doi.org/10.1080/01621459.2019.1677241>
- 582 Collins, L. M., Schafer, J., & Kam, C. (2001). A Comparison of Restrictive Strategies in
583 Modern Missing Data Procedures. *Psychological Methods*, 6(June).
584 <https://doi.org/10.1037/1082-989X.6.4.330>
- 585 Dornelas, M., Gotelli, N. J., McGill, B., Shimadzu, H., Moyes, F., Sievers, C., & Magurran,
586 A. E. (2014). Assemblage time series reveal biodiversity change but not systematic loss.

- 587 *Science*, 344(6181), 296–299. <https://doi.org/10.1126/science.1248484>
- 588 Elliott, M. R., & Valliant, R. (2017). Inference for nonprobability samples. *Statistical*
589 *Science*, 32(2), 249–264. <https://doi.org/10.1214/16-STS598>
- 590 Forister, M. L., Black, S. H., Elphick, C. S., Grames, E. M., Halsch, C. A., Schultz, C. B., &
591 Wagner, D. L. (2023). Missing the bigger picture: Why insect monitoring programs are
592 limited in their ability to document the effects of habitat loss. *Conservation Letters*,
593 *September 2022*, 1–6. <https://doi.org/10.1111/conl.12951>
- 594 Geldmann, J., Heilmann-Clausen, J., Holm, T. E., Levinsky, I., Markussen, B., Olsen, K.,
595 Rahbek, C., & Tøttrup, A. P. (2016). What determines spatial bias in citizen science?
596 Exploring four recording schemes with different proficiency requirements. *Diversity and*
597 *Distributions*, 22(11), 1139–1149. <https://doi.org/10.1111/ddi.12477>
- 598 Gelman, A. (2007). Struggles with survey weighting and regression modeling. *Statistical*
599 *Science*, 22(2), 153–164. <https://doi.org/10.1214/088342306000000691>
- 600 Gelman, A., & Little, T. (1997). poststratification into many categories using hierarchical
601 regression. *Survey Methodology*, 23(2), 127–335.
- 602 Girardello, M., Chapman, A., Dennis, R., Kaila, L., Borges, P. A. V., & Santangeli, A.
603 (2019). Gaps in butterfly inventory data: A global analysis. *Biological Conservation*,
604 236(November 2018), 289–295. <https://doi.org/10.1016/j.biocon.2019.05.053>
- 605 Gonzalez, A., Cardinale, B. J., Allington, G. R. H., Byrnes, J., Endsley, K. A., Brown, D. G.,
606 Hooper, D. U., Isbell, F., O'Connor, M. I., & Loreau, M. (2016). Estimating local
607 biodiversity change: A critique of papers claiming no net loss of local diversity.
608 *Ecology*, 97(8), 1949–1960. <https://doi.org/10.1890/15-1759.1>
- 609 Gregory, R. D., Van Strien, A., Vorisek, P., Meyling, A. W. G., Noble, D. G., Foppen, R. P.
610 B., & Gibbons, D. W. (2005). Developing indicators for European birds. *Philosophical*
611 *Transactions of the Royal Society B: Biological Sciences*, 360(1454), 269–288.
612 <https://doi.org/10.1098/rstb.2004.1602>
- 613 Hudson, L. N., Newbold, T., Contu, S., Hill, S. L. L., Lysenko, I., Palma, D., Phillips, H. R.
614 P., Senior, R. A., Bennett, D. J., Booth, H., Garon, M., Michelle, L., Correia, D. L. P.,
615 Day, J., Echeverr, S., Harrison, K., Ingram, D. J., Jung, M., Kemp, V., ... Fernando, A.
616 B. (2014). *The PREDICTS database : a global database of how local terrestrial*
617 *biodiversity responds to human impacts*. 4701–4735.
- 618 Hughes, A., Orr, M., Ma, K., Costello, M., Waller, J., Provoost, P., Zhu, C., & Qiao, H.
619 (2020). Sampling biases shape our view of the natural world. *Ecography*, 44, 1259–
620 1269. <https://doi.org/10.1111/ecog.05926>
- 621 Intermap. (2009). *NEXTMap British Digital Terrain 50m resolution (DTM10) Model Data by*
622 *Intermap*. NERC Earth Observation Centre.
623 <https://catalogue.ceda.ac.uk/uuid/f5d41db1170f41819497d15dd8052ad2>
- 624 Johnston, A., Matechou, E., & Dennis, E. B. (2022). Outstanding challenges and future
625 directions for biodiversity monitoring using citizen science data. *Methods in Ecology*
626 *and Evolution*, *February*. <https://doi.org/10.1111/2041-210X.13834>
- 627 Johnston, A., Moran, N., Musgrove, A., Fink, D., & Baillie, S. R. (2020). Estimating species
628 distributions from spatially biased citizen science data. *Ecological Modelling*,
629 422(December 2019), 108927. <https://doi.org/10.1016/j.ecolmodel.2019.108927>

- 630 Lohr, S. (2022). *Sampling: Design and analysis* (3rd ed.). CRC Press.
- 631 Lumley, T. (2010). *Complex surveys: A Guide to Analysis Using R* (1st ed.). Wiley.
- 632 Makela, S., Si, Y., & Gelman, A. (2014). Statistical Graphics for Survey Weights. *Revista*
633 *Colombiana de Estadística*, 37(2Spe), 285–295.
634 <https://doi.org/10.15446/rce.v37n2spe.47937>
- 635 Mandeville, C. P., Nilsen, E. B., & Finstad, A. G. (2022). Spatial distribution of biodiversity
636 citizen science in a natural area depends on area accessibility and differs from other
637 recreational area use. *Ecological Solutions and Evidence*, 3(4), 1–14.
638 <https://doi.org/10.1002/2688-8319.12185>
- 639 McClure, C. J. W., & Rolek, B. W. (2023). Pitfalls arising from site selection bias in
640 population monitoring defy simple heuristics. *Methods in Ecology and Evolution*, 14(6),
641 1489–1499. <https://doi.org/10.1111/2041-210X.14120>
- 642 Meineke, E. K., & Daru, B. H. (2021). Bias assessments to expand research harnessing
643 biological collections. *Trends in Ecology and Evolution*, 36(12), 1071–1082.
644 <https://doi.org/10.1016/j.tree.2021.08.003>
- 645 Meng, X.-L. (2018). Statistical paradises and paradoxes in big data (I): Law of large
646 populations, big data paradox, and the 2016 us presidential election. *Annals of Applied*
647 *Statistics*, 12(2), 685–726. <https://doi.org/10.1214/18-AOAS1161SF>
- 648 Meng, X.-L. (2022). Comments on the Wu (2022) paper by Xiao-Li Meng 1 : Miniaturizing
649 data defect correlation : A versatile strategy for handling non-probability samples.
650 *Survey Methodology*, 48(2), 1–22.
- 651 Mercer, A., Lau, A., & Kennedy, C. (2018). For Weighting Online Opt-In Samples, What
652 Matters Most? In *Pew Research Center* (pp. 1–55). <https://pewrsr.ch/3heqknn>
- 653 Morton, R., Marston, C., O’Neil, A., & Rowland, C. (2022). *Land Cover Map 2018 (1km*
654 *summary rasters, GB and N. Ireland)*. NERC EDS Environmental Information Data
655 Centre. <https://doi.org/https://doi.org/10.5285/9b68ee52-8a95-41eb-8ef1-8d29e2570b00>
- 656 Nelson, G., & Ellis, S. (2019). The history and impact of digitization and digital data
657 mobilization on biodiversity research. *Philosophical Transactions of the Royal Society*
658 *B: Biological Sciences*, 374(1763), 2–10. <https://doi.org/10.1098/rstb.2017.0391>
- 659 Outhwaite, C., Gregory, R. D., Chandler, R. E., Collen, B., & Isaac, N. J. B. (2020). Complex
660 long-term biodiversity change among invertebrates, bryophytes and lichens. *Nature*
661 *Ecology & Evolution*. <https://doi.org/10.1038/s41559-020-1111-z>
- 662 Perring, F., & Walters, S. (1962). *Atlas of the British flora*. Thomas Nelson & sons.
- 663 Pescott, O. L., Stroh, P. A., Humphrey, T. A., & Walker, K. J. (2022). Simple methods for
664 improving the communication of uncertainty in species ’ temporal trends. *Ecological*
665 *Indicators*, 141(May). <https://doi.org/https://doi.org/10.1016/j.ecolind.2022.109117>
- 666 Pescott, O. L., Walker, K. J., Harris, F., New, H., Cheffings, C. M., Newton, N., Jitlal, M.,
667 Redhead, J., Smart, S. M., & Roy, D. B. (2019). The design, launch and assessment of a
668 new volunteer-based plant monitoring scheme for the United Kingdom. *PLoS ONE*,
669 14(4), 1–30. <https://doi.org/10.1371/journal.pone.0215891>
- 670 Pescott, O. L., Walker, K. J., Pocock, M. J. O., Jitlal, M., Outhwaite, C. L., Cheffings, C. M.,

- 671 Harris, F., & Roy, D. B. (2015). Ecological monitoring with citizen science: The design
672 and implementation of schemes for recording plants in Britain and Ireland. *Biological*
673 *Journal of the Linnean Society*, 115(3), 505–521. <https://doi.org/10.1111/bij.12581>
- 674 Powney, G. D., Carvell, C., Edwards, M., Morris, R. K. A., Roy, H. E., Woodcock, B. A., &
675 Isaac, N. J. B. (2019). Widespread losses of pollinating insects in Britain. *Nature*
676 *Communications*, 10(2019), 1–6. <https://doi.org/10.1038/s41467-019-08974-9>
- 677 Preston, C.D., Pearman, D.A. & Dines, T. D. (2002). *New Atlas of the British and Irish*
678 *Flora*. (eds). Oxford University Press.
- 679 Rowland, C., Marston, C., Morton, R., & O’Neil, A. (2020). *Land Cover Map 1990 (1km*
680 *dominant target class, GB) v2*. NERC EDS Environmental Information Data Centre.
681 <https://doi.org/https://doi.org/10.5285/f5e3bd00-efd0-4dc6-a454-aa597d84764a>
- 682 Rubin, D. B. (1976). Inference and missing data. *Biometrika*, 63(3), 581–592.
683 <https://doi.org/https://doi.org/10.1093/biomet/63.3.581>
- 684 Ruete, A. (2015). Displaying bias in sampling effort of data accessed from biodiversity
685 databases using ignorance maps. *Biodiversity Data Journal*, 3(1), 1–15.
686 <https://doi.org/10.3897/BDJ.3.e5361>
- 687 Stroh, P. A., Walker, K., Humphrey, T. A., Pescott, O. L., & Burkmar, R. (2023). *Plant Atlas*
688 *2020: Mapping Changes in the Distribution of the British and Irish Flora*. Princeton
689 Univ. Press.
- 690 Tchetgen Tchetgen, E. J., & Wirth, K. E. (2017). A general instrumental variable framework
691 for regression analysis with outcome missing not at random. *Biometrics*, 73(4), 1123–
692 1131. <https://doi.org/10.1111/biom.12670>
- 693 Thoemmes, F., & Mohan, K. (2015). Graphical Representation of Missing Data Problems.
694 *Structural Equation Modeling*, 22(4), 631–642.
695 <https://doi.org/10.1080/10705511.2014.937378>
- 696 Thoemmes, F., & Rose, N. (2014). A Cautious Note on Auxiliary Variables That Can
697 Increase Bias in Missing Data Problems. *Multivariate Behavioral Research*, 49(5), 443–
698 459. <https://doi.org/10.1080/00273171.2014.931799>
- 699 UNEP-WCMC, & IUCN. (2020). *Protected Planet: The World Database on Protected Areas*
700 *(WDPA)/The Global Database on Protected Areas Management Effectiveness*.
701 <https://www.protectedplanet.net/en/thematic-areas/wdpa>
- 702 Valliant, R. (2020). Comparing Alternatives for Estimation from Nonprobability Samples.
703 *Journal of Survey Statistics and Methodology*, 8(2), 231–263.
704 <https://doi.org/10.1093/jssam/smz003>
- 705 Valliant, R., Dever, J. A., & Kreuter, F. (2018). *Practical tools for designing and weighting*
706 *survey samples* (2nd ed.). Springer Cham. [https://doi.org/https://doi.org/10.1007/978-3-](https://doi.org/https://doi.org/10.1007/978-3-319-93632-1)
707 [319-93632-1](https://doi.org/https://doi.org/10.1007/978-3-319-93632-1)
- 708 van Strien, A. J., & van Grunsven, R. H. A. (2023). In the past 100 years dragonflies declined
709 and recovered by habitat restoration and climate change. *Biological Conservation*,
710 277(December 2022), 109865. <https://doi.org/10.1016/j.biocon.2022.109865>
- 711 Van Swaay, C. A. M., Nowicki, P., Settele, J., & Van Strien, A. J. (2008). Butterfly
712 monitoring in Europe: Methods, applications and perspectives. *Biodiversity and*

- 713 *Conservation*, 17(14), 3455–3469. <https://doi.org/10.1007/s10531-008-9491-4>
- 714 Van Swaay, C. A. M., Plate, C. L., & Van Strien, A. J. (2002). Monitoring butterflies in the
715 Netherlands: how to get unbiased indices. *Proceedings of the Section Experimental and*
716 *Applied Entomology of the Netherlands Entomological Society*, 13, 21–27.
- 717 Vellend, M., Baeten, L., Myers-Smith, I. H., Elmendorf, S. C., Beauséjour, R., Brown, C. D.,
718 De Frenne, P., Verheyen, K., & Wipf, S. (2013). Global meta-analysis reveals no net
719 change in local-scale plant biodiversity over time. *Proceedings of the National Academy*
720 *of Sciences of the United States of America*, 110(48), 19456–19459.
721 <https://doi.org/10.1073/pnas.1312779110>
- 722 Weiser, E. L., Diffendorfer, J. E., Lopez-Hoffman, L., Semmens, D., & Thogmartin, W. E.
723 (2020). Challenges for leveraging citizen science to support statistically robust
724 monitoring programs. *Biological Conservation*, 242(October 2019).
725 <https://doi.org/10.1016/j.biocon.2020.108411>
- 726 Wu, C. (2022). *Statistical inference with non-probability survey samples*. 12.
- 727 Wu, C., & Sitter, R. R. (2001). A model-calibration approach to using complete auxiliary
728 information from survey data. *Journal of the American Statistical Association*, 96(453),
729 185–193. <https://doi.org/10.1198/016214501750333054>
- 730