# Descriptive inference using large, unrepresentative nonprobability samples: An introduction for ecologists

[*1]Robin J. Boyd, [2]Gavin B. Stewart and [1]Oliver L. Pescott

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

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

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

## Abstract

In the age of big data, it is essential to remember that the size of a dataset is not all that matters. This is particularly true where the goal is to draw inferences about some wider population, in which case it is far more important that the data are *representative* of that population. It is possible to adjust unrepresentative samples so that they more closely resemble the population in terms of "auxiliary variables". If the auxiliaries predict sample inclusion and/or the variable of interest well, then the adjusted sample estimates will be closer to the truth. Several survey sampling techniques exist to perform such adjustments, but most are not familiar to ecologists. We applied five types of adjustment—subsampling, quasi-randomisation, poststratification, superpopulation modelling, and multilevel regression and poststratification—to a simple two-part biodiversity monitoring problem. The first part was to estimate mean occupancy of the plant *Calluna vulgaris* in Great Britain in two time-periods (1987-1999 and 2010-2019); the second was to estimate the difference between the two (i.e. the trend). *Calluna vulgaris* is an attractive case study because we have good estimates of its true distribution in both time-periods. We estimated the means and trend using large, but (originally) unrepresentative, samples. Compared to the unadjusted estimates, the means and trends estimated using most adjustment methods were more accurate, although their uncertainty intervals generally did not cover the true values. Quasi-randomisation performed especially poorly, and we explain why. Most adjustments were far more successful at bringing the distributions of the auxiliary variables in the samples closer to those in the population than they were at improving the estimates of population means and trends. This implies that the major challenge for adjusting unrepresentative samples in biodiversity monitoring is assembling a suitable set of auxiliary variables (i.e. predictors of sample inclusion and the variable of interest). This challenge will be particularly acute for poorly studied taxa and those whose habitat requirements or sampling biases are not reflected in available data.

## Introduction

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

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

45  More important than the size of a sample for descriptive inference is whether it is representative of the
46  population (X. L. Meng, 2018). In a representative sample, the distribution of the variable of interest
47  is similar to its distribution in the population (Bethlehem et al., 2008). An equivalent definition is that
48  there is little to no correlation between inclusion in the sample and the variable of interest—the "data
49  defect correlation" (ddc; Meng, 2018). Intuitively, statistics derived from a representative sample,
50  such as means and proportions, will be similar to their population equivalents. The challenge is that
51  the variable of interest is unknown in non-sampled population units, so it is typically impossible to
52  calculate a sample's representativeness exactly.

53  Rather than measuring sample representativeness in terms of the variable of interest, which is not
54  known for all population units, it is possible to approximate it using "auxiliary variables". Auxiliary
55  variables are those that are thought to predict the variable of interest or the probability that each
56  population unit was sampled. Such variables might be available for every population unit. For
57  example, climate variables might explain a species' occupancy, and data on these variables are
58  available for every 1 km$^2$ grid square across the globe (Fick & Hijmans, 2017). If the distributions of
59  auxiliary variables in the sample are different to those in the population, then the sample is likely to be
60  unrepresentative, at least with respect to those variables (Bethlehem et al., 2008).

61  It is possible to adjust an unrepresentative sample by placing more weight on population units that
62  were less likely to be sampled. Weighting is simple where the probability that each population unit
63  was sampled is known (i.e. in a probability sample). For example, rather than using a sample mean to
64  estimate a population mean, the researcher would use a weighted mean with the weights being the
65  inverses of the inclusion probabilities (Lohr, 2022). Weighting of this type is known as "design-based
66  inference", because the inclusion probabilities are a feature of the sampling design. Unfortunately,
67  design-based inference is not applicable for many "big" biodiversity datasets, whose sample inclusion
68  probabilities are not known (i.e. they are nonprobability samples), so alternatives are required.

69  Most approaches to descriptive inference from nonprobability samples make use of auxiliary
70  variables. The details differ, but the general strategy is to weight the sample in such a way that the
71  distributions of the auxiliary variables in the sample more closely resemble those in the population
72  (Valliant et al., 2018). If the auxiliaries predict the variable of interest and sample inclusion well, then
73  this is essentially the same as bringing the distribution of the variable of interest in the sample closer
74  to its distribution in the population (i.e. making the sample more representative; see X. Meng, [2022],
75  who demonstrates this mathematically).

76  Use of sample adjustments in biodiversity monitoring is variable. It is common for monitoring
77  schemes to weight samples in such a way that the relative frequencies of habitats or geographic areas
78  in the sample are similar to those in the population (Gregory et al., 2005; C.A.M. Van Swaay et al.,
79  2002; Chris A.M. Van Swaay et al., 2008; Weiser et al., 2020). But it is also common to see sample
80  representativeness ignored, an issue that has led to some high-profile controversies in the biodiversity
81  monitoring literature (Boyd, Powney, et al., 2023). We suspect that many of those who do not deal
82  with issues of sample representativeness are not familiar with the gravity of the problem or the
83  relevant theory and adjustment methods.

84  In this paper, we introduce five approaches to descriptive inference using unrepresentative
85  nonprobability samples and demonstrate how they relate to each other (conceptually and
86  mathematically). We apply each approach to a simple two-part biodiversity monitoring problem. The
87  first part is to estimate mean occupancy of the plant *C. vulgaris* across 1 km grid squares in Britain in
88  two time-periods; the second is to estimate the difference between the two (i.e. the trend). *Calluna*
89  *vulgaris* is an attractive case study because we have good estimates of its true geographic distribution
90  in both periods from several sources. The approaches to inference that we demonstrate are
91  subsampling, quasi-randomisation (Elliott and Valliant, 2017), poststratification (Little, 1993),
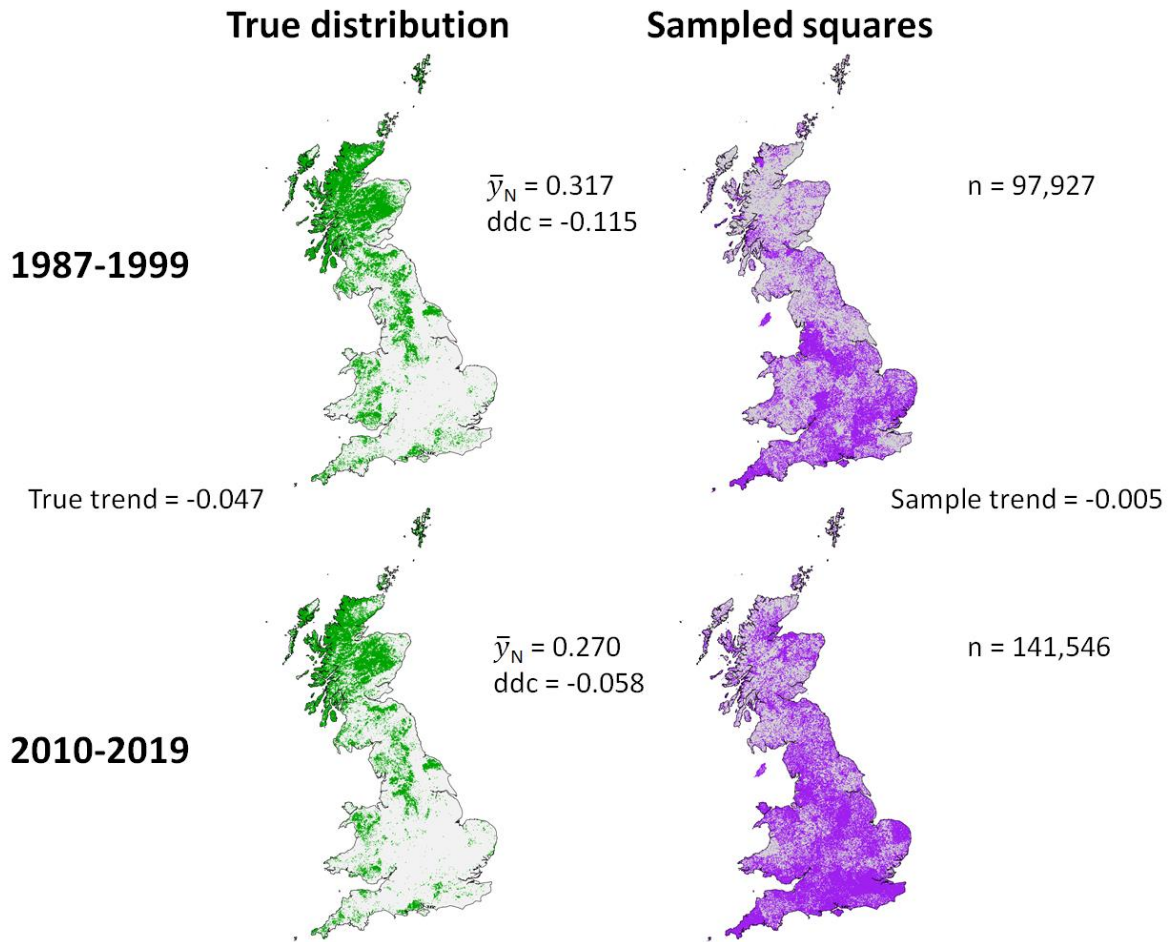
92  superpopulation modelling (Valliant, 2009) and Multilevel Regression and Poststratification (MRP;
93  Gelman, 2007; Gelman and Little, 1997). Each can be (MRP more loosely than the rest) interpreted as
94  an attempt to weight the sample in such a way that it more closely resembles the population, in the
95  hope that this results in more accurate descriptive inferences. We demonstrate the effects of each
96  approach on the distributions of auxiliary variables in the sample, as well as on the resulting estimates
97  of mean occupancy in each period and the time trend between the two. Applying the adjustment
98  methods to a real-world example reveals challenges that ecologists are likely to face, and we discuss
99  these in detail.

# Methods

## Estimating the true distribution of *Calluna vulgaris*

102  We approximated the true distribution of the dwarf shrub vascular plant *Calluna vulgaris* (Heather) in
103  two time periods: 1987–1999 and 2010–2019. For the first period, we used the 1990 UKCEH land
104  cover map (Rowland et al., 2020); for the second, we used the 2018 version (Morton et al., 2022).
105  From these maps, we identified 1 km grid squares (British National Grid, EPSG:27700) with >0%
106  heather or heather grassland cover. To these, we added 1 km squares in which *C. vulgaris* was
107  recorded in each time period by the Botanical Society of Britain and Ireland (BSBI); the time periods
108  used cover the main periods of recording for two national distribution atlases (Preston, C.D., Pearman,
109  D.A. & Dines, 2002; Stroh et al., 2023). Acknowledging that some 1 km squares may have been
110  erroneously classed as having some heather or heather grassland coverage by the land cover maps, we
111  removed any 1 km squares that fell within 10 km grid squares in which *C. vulgaris* had not been
112  recorded by the BSBI in the period 1950–2019. Given that this period includes recording for three
113  national distribution atlases (the two cited above plus Perring & Walters, 1962), we assume that the
114  union of all 10 km occurrences within this period encompasses all known populations irrespective of
115  finer scale changes. Figure 1 maps the resulting estimates of the true 1 km distributions of *C. vulgaris*
116  in both time-periods.

117

True distribution　　　Sampled squares

$\bar{y}_N = 0.317$
ddc = -0.115

n = 97,927

1987-1999

True trend = -0.047

Sample trend = -0.005

$\bar{y}_N = 0.270$
ddc = -0.058

n = 141,546

2010-2019

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

## Sample data on *Calluna vulgaris* occupancy

The 1 km samples for both time periods ("Sampled squares in Fig. 1) encompass any vascular plant data assigned to a single day, either at the 1 km scale or finer, collected by the BSBI for the national distribution atlases of Preston et al. (2002) and Stroh et al. (2023).

## Auxiliary data

We used five auxiliary variables for which data are available for all 1 km grid squares in Great Britain (Table 1). Most of the auxiliary variables indicate the accessibility or attractiveness of grid squares, which tend to be associated with site selection in citizen science datasets (Geldmann et al., 2016). Elevation is a potential predictor of *C. vulgaris* occupancy.

Originally, we included three additional predictors of *C. vulgaris* occupancy—the first and third principal components of climate space in Britain and soil pH—but later omitted them. We had previously found the climate variables to be important predictors of 1 km habitat suitability for *C. vulgaris* using species distribution models (Boyd, Harvey, et al., 2023). Including these predictors did not improve the estimates of mean occupancy, a point that we expand on in the discussion. Reducing

141  the set of auxiliary variables simplifies matters for some of the adjustment methods that we present
142  below.

143  For simplicity, we assume that the auxiliary variables are constant between time-periods. This
144  assumption is obviously violated for some variables (e.g. road length and postcode density). However,
145  this should not matter if, in reality, the variables in period one are correlated with those in period two,
146  because any given grid square will generally have a higher or lower value than the others regardless of
147  the period. We think that this situation is plausible: for example, there is a higher density of postcodes
148  in London in period two than in period one, but in either period, it has a higher density than
149  elsewhere. Another reason to use one set of auxiliary data for both time-periods is to make our
150  findings more applicable to situations in which temporally resolved data are not available (e.g. in data
151  poor countries or periods in the distant past).

152  Table 1. Auxiliary variables used for sample adjustment.

| Variable | Reason for inclusion | Details | Reference(s) |
|---|---|---|---|
| Postcode density | Indicates population density in vicinity | Total number of postcodes in the focal grid square and its 299 nearest neighbours | ONS (2021) |
| Road length | Indicates accessibility | The total length of all "Roads" and "Link roads" ("Highways" class of the OpenStreetMap ontology) in the focal grid square and its 299 nearest neighbours | Data from https://www.openstreetmap.org/ under an open database license |
| Proportion in protected area | Indicates potential attractiveness to surveyor | Proportion of the focal grid square with some level of "protection". Includes everything from SSSIs to e.g. local nature reserves | UNEP-WCMC & IUCN (2020) |
| Proportion open access land | Indicates accessibility | Proportion of land legally designated as open access within 1 km grid square | All open access land datasets in GB are available via an Open Government License. For England, we used the CRoW act 2000 layer. For Wales, we combined the registered common land, other statutory access land, open country and public forest datasets. All of the Scottish countryside is open access. |
| Average elevation | Predictor of *C. vulgaris* occupancy | Average elevation of 1 km grid square calculated from 50 m digital terrain model | Intermap (2009) |

153

## Estimating the per-period population mean

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

$$\bar{y}_N = \frac{1}{N}\sum_{i=1}^{N} y_i = \sum_{i=1}^{N}\frac{y_i}{N} = \sum_{i=1}^{N}\frac{y_i \, w_i}{\sum_N w_i}, \qquad\qquad 1)$$

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

## The design-based estimator

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

$$\bar{y}_{db} = \sum_{i=1}^{n}\frac{y_i \, w_i}{\sum_n w_i}. \qquad\qquad 2)$$

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

Sample inclusion probabilities are, by definition, not known for nonprobability samples, so alternative estimators are required. We present five such estimators below, three of which–quasi-randomisation, poststratification and superpopulation modelling–are explicit attempts to come up with a set of weights $w_i$ that produce a reasonable estimate of $\bar{y}_N$ under 2). The other two, subsampling and MRP, are conceptually similar.

## Estimators for nonprobability samples

The following estimators are used in survey sampling to estimate population means from nonprobability samples. More detail on each can be found in Valliant et al. (2018), Lumley (2010) and Lohr (2022).

### *Naïve sample mean*

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

### *Quasi-randomisation*
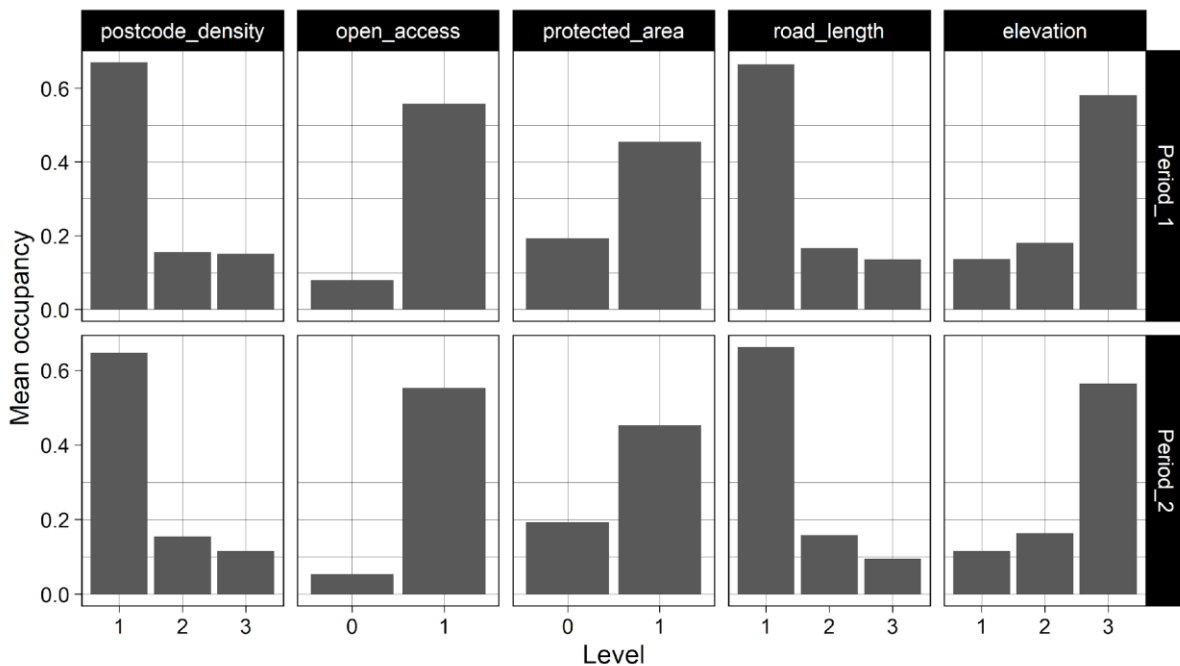
An alternative approach is to imagine that the nonprobability sample was selected probabilistically and to estimate the implied inclusion probabilities. Any binary model and auxiliary data can be used. Once inclusion probabilities $p_i$ have been estimated, the weights $w_i = 1/p_i$ (as in the design-based estimator). In our example, we used random forests and the auxiliary data in Table 1 to estimate pseudo-inclusion probabilities. More complex appraoches are possible and have been used to map species distributions (Johnston et al., 2020).

*Poststratification*
194    Another approach to estimating sampling weights is poststratification. Poststratification requires
195    categorical auxiliary data, so continuous variables must be discretized prior to analysis (Valliant,
196    2020). The auxiliary variables are crossed (think contingency tables) to create poststrata. Each
197    poststratum $j$ has a sample size $n_j$ and population size $N_j$. The sampling weight $w_i$ for population unit
198    $i$ in poststratum $j$ is given by $N_j/n_j$.

199    In our example, we split most auxiliary variables into three categories using their terciles (i.e. cut
200    points at the 33rd and 67th percentiles). This did not make sense for the variables denoting the
201    proportion of each grid square that is open access land and protected area, because most squares took
202    the value one or zero. We split these variables into two categories, 0 and >0, i.e. whether or not there
203    is some open or protected land in the grid square. Discretization initially gave $3 \times 3 \times 3 \times 2 \times$
204     $2 = 108$ poststrata, from which we subtracted one poststratum that contained no population units,
205    leaving 107.

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



210

211    Figure 2. Mean occupancy of *Calluna vulgaris* for each level of the auxiliary variables (Table 1) in
212    each time-period. The auxiliary variables were originally on a continuous scale, but we discretized
213    them to enable poststratification. See the main text for details.

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

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

$$E_M(y_i) = \boldsymbol{x}_i^T \beta, \qquad\qquad 3)$$

221 where the subscript $M$ indicates that the expectation (mean) is with respect to the model, $\boldsymbol{x}_i$ is a vector
222 of predictors for unit $i$, the superscript $T$ indicates that the vector $\boldsymbol{x}_i$ has been transposed (to a row
223 vector) and $\beta$ is a column vector of parameters. There is some matrix notation in 3) and what follows,
224 but the logic should be apparent to those who do not understand the precise detail. A prediction of $y$
225 for unit $i$ is

$$\hat{y}_i = x_i^T \hat{\beta}. \qquad\qquad 4)$$

226 The accent on β indicates that it is an estimate. Given a sample s, one estimator of $\beta$ is $\hat{\beta} = A_s^{-1} X_s^T y_s$,
227 where $A_s = X_s^T X_s$, $X_s$ is an $n\ x\ p$ matrix of covariates and $\boldsymbol{y}_s$ is an $n$ vector of $y$'s (Valliant, 2020). If
228 $\bar{s}$ is the set of non-sampled population units, the superpopulation model prediction of the population
229 mean is

$$\bar{y}_{sp} = \frac{\sum_{i \in s} y + \sum_{i \in \bar{s}} \hat{y}}{N}. \qquad\qquad 5)$$

230 That is, it is the sum of the known outcome values in the sample and those predicted by the model for
231 the remainder of the population divided by the population total. A feature of $\bar{y}_{sp}$ is that it can be
232 expressed in the same form as the design-based estimator in 2), with $w_i = 1 + \boldsymbol{t}_{\bar{s}}^x A_s^{-1} x_i$ and $\boldsymbol{t}_{\bar{s}}^x$ being
233 the vector of population totals of the auxiliary variables in non-sampled population units (Elliott and
234 Valliant, 2017). (Code to verify this numerically is available at
235 https://github.com/robboyd/selectionBiasEffects/tree/master/R.) Like the other adjustment models,
236 then, the superpopulation estimator is an approach to estimating the sampling weights $w_i$.

237 Linear regression might seem like an unusual choice of model for a binary outcome (occupancy), but
238 we felt that it was the best option here. One reason is that the implied model is actually linear for an
239 estimator of the form 2) (Valliant, 2020). Most important, however, is that the use of a linear model
240 enables the estimation of sampling weights (Valliant et al., 2018; supplementary material 1). This is
241 helpful, because those weights can be used to show the effects of superpopulation modelling on the
242 distributions of the auxiliary variables in the sample (see "Evaluating the effects of the adjustments"
243 below). Alternative models can be used where weights are not required (e.g. Wu and Sitter, 2001). In
244 our example superpopulation model, we used the auxiliary variables in Table 1 as predictors.

245 *Subsampling*
246 Perhaps more familiar to ecologists than the above approaches is subsampling (Beck et al., 2014;
247 Steen et al., 2020). The idea is to create a representative "miniature" of the population out of the
248 sample (Meng, 2022) and to calculate the quantity of interest (mean occupancy) from this subsample.
249 Subsampling trades sample size for representativeness.

250 Our approach was to draw weighted random samples of size 500 with replacement from the original
251 samples [note that these weights are different to sampling weights in 2)]. The decision to set $n = 500$
252 was somewhat arbitrary, but changing the subsample size makes little difference to the point estimates
253 of the population means (although they become more precise with increasing subsample size;
254 supplementary material 1). We assigned each grid square $i$ in poststratum $j$ (using the same strata as
255 described above under Poststratification) a weight of $n_j/N_j$. The result was subsamples whose
256 members were more likely to be from strata comprising a larger fraction of the population. The
257 subsample mean is the estimator of the population mean.

258 Rather than using a single subsample, we repeated the process 1000 times and used the mean of the
259 estimated means (i.e. bootstrapping). This was necessary, because the estimated means were sensitive
260 to the random component of the subsampling.

261 *Multilevel regression and poststratification (MRP)*
262 MRP is an extension of poststratification and a variation of superpopulation modelling (Gelman,
263 2007; Gelman & Little, 1997; Valliant et al., 2018). A hierarchical model is used to estimate mean
264 occupancy in each poststratum. The advantage of using a hierarchical model is that cells with few or
265 no data borrow information from cells with more data (i.e. partial pooling or shrinkage is exploited).
266 The population mean is the weighted mean of the stratum means, where the weights are equal to the
267 proportion of the population in each stratum.

268 Our hierarchical model is a simple one. It is a binomial GLM with a logit link function, a fixed
269 intercept and a random intercept for each auxiliary variable (see https://mc-
270 stan.org/rstanarm/articles/mrp.html for a similar formulation). A more complex model might include
271 interactions among the auxiliaries (e.g. Ghitza and Gelman, 2013), but we found these take several
272 times longer to run. Long run times may be undesirable for production-type statistical workflows in
273 biodiversity monitoring, where models might need to be fitted for thousands of species in tens of
274 time-periods. Even without interactions, and on a computer cluster, the models took around ten hours
275 to run per time-period. We fitted the model in a Bayesian framework using 5 Markov Chain Monte
276 Carlo (MCMC) chains, each with 1000 iterations. This was sufficient to achieve convergence on all
277 parameters in both time-periods.

278 *Confidence intervals*
279 We present 95% confidence/credible intervals for all estimates of mean occupancy (credible intervals
280 for MRP, which was implemented in a Bayesian framework). The *survey* package (Lumley, 2010),
281 which we used to calculate the sample means, the superpopulation model estimates and the
282 poststratified estimates, calculates the confidence intervals automatically. It accounts for the sampling
283 weights where relevant. We used percentile confidence intervals from the bootstrapped subsamples.

284 ## Estimating the trend in mean occupancy
285 Having estimated mean occupancy in each time-period, the next step was to estimate the difference
286 between the two $= \bar{y}_2 - \bar{y}_1$ (i.e. the trend). The standard errors of the trends are

287 $\sqrt{var\left(\underline{y}_2\right) + var\left(\underline{y}_1\right)}$ (Gelman, 2007), where the variances are sampling *not sample* variances (i.e.

288 the square of the standard error rather than a measure of variability in the samples). We used the
289 standard errors returned by the *survey* package, which accounts for the sampling weights. We present
290 95% confidence intervals for the trends from most estimators ($\pm 1.96 \times$ the standard errors). MRP is
291 one exception, because the 95% credible interval can be calculated directly from the posterior
292 distribution of $\bar{y}_2 - \bar{y}_1$. Similarly, we extracted percentile 95% confidence intervals for the
293 subsampling estimator from the bootstrapped distribution of trends.

294 ## Evaluating the effects of the adjustments
295 We used relative frequency plots (c.f. Makela et al., 2014) to assess whether the adjustments brought
296 the distributions of the auxiliary variables in the samples closer to their distributions in the population.
297 The first step was to split each auxiliary variable into fifty bins of equal width spanning its range. The
298 relative frequency of grid squares (the $i$'s) in each bin k is $N_{i,k}/N$, where $N_{i,k}$ is the number of grid
299 squares in each bin $k$ in the population and $N$ is the population size (we use $k$ to index the bins to
300 distinguish them from the strata described earlier). Similarly, the relative frequency of sampled grid
301 squares in each $k$ is $n_{i,k}/n$, where $n_{i,k}$ is the number of sampled grid squares in bin $k$ and $n$ is the
302 total sample size. In the adjusted samples, the equivalent relative frequency is $\frac{\sum_{i \in k} w_i}{\sum_N w_i}$ (slightly
303 different for subsampling; see below). We compared the original and adjusted samples' deviations
304 from the population using the Mean Absolute Error (MAE) of the relative frequencies across all $k$. If
305 the MAE from the adjusted sample is smaller than the original sample, then the adjustment brought
306 the distribution of the auxiliary variable closer to its population distribution.

307    We were not able to construct adjusted relative frequency distributions from MRP so omit it from this
308    portion of the analysis. The problem is that, whilst it has been shown how to derive unit-level
309    sampling weights where the multilevel model is linear (Gelman, 2007), no formula has yet been
310    derived for the case of the binomial GLM (Valliant et al., 2018). There is no obvious way to derive
311    weights from the subsampling estimator either. However, for this estimator, the adjusted relative
312    frequencies of the auxiliaries are simply their distributions in the subsamples so are simple to obtain.

313    Assessing whether the estimates of mean occupancy in each period and the trend were improved by
314    each adjustment method was simpler. We measured the difference between the point estimates of
315    mean occupancy and the truth using the absolute error = $|\bar{y}_N - \bar{y}_{est}|$, where $\bar{y}_{est}$ is the estimate. For
316    the trends, whose signs are of interest, we simply used the differences between the estimates and the
317    truth. We also assessed whether the confidence/credible intervals produced by each method covered
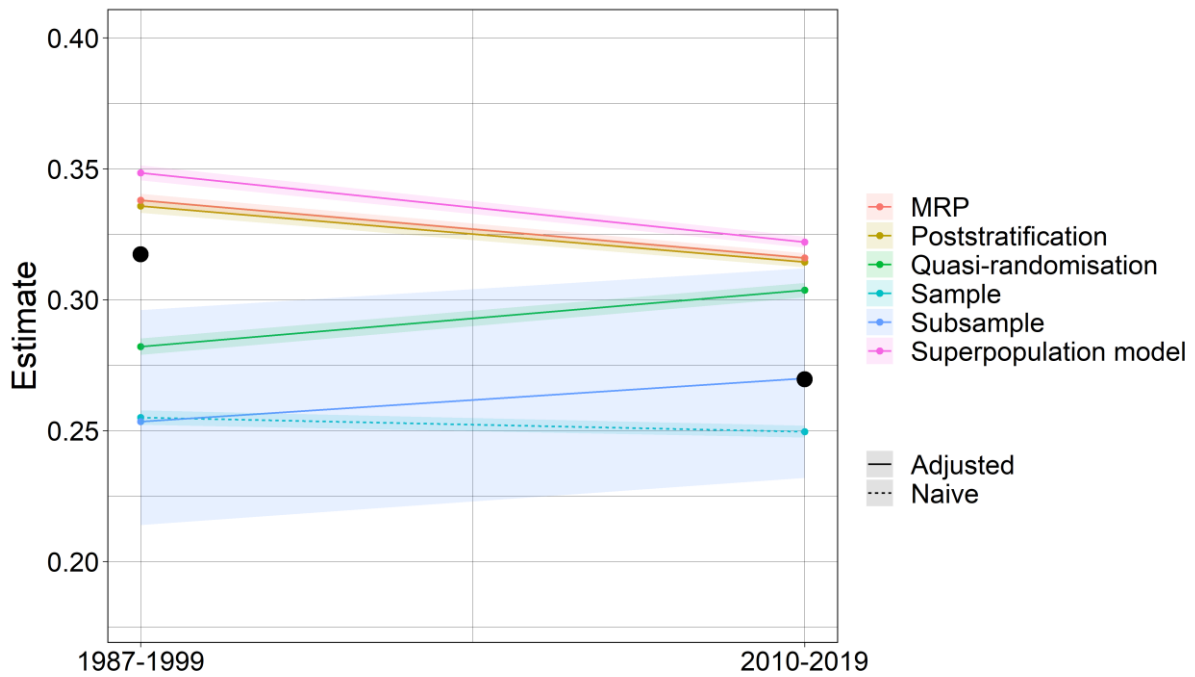318    the true means and trend.

319    # Results

320    ## Per-period sample representativeness and estimated mean occupancy
321    The samples are large but somewhat unrepresentative (Fig. 1). Forty-three percent of grid squares
322    were sampled in period one, and the ddc is -0.115; in period two, 62% of grid squares were sampled,
323    and the ddc is -0.058. A consequence of these ddc's is that the naive sample means underestimate the
324    population means, especially in period one where the magnitude of the ddc is greater (Fig. 3).

325    The adjustment methods did not always result in improved point estimates of mean occupancy
326    relative to the naive sample means (Fig. 3). In period one, the adjusted estimates were generally better
327    in terms of absolute errors, with the exception of the subsample estimate, which was worse. In period
328    two, on the other hand, the estimate from the subsample was the only one to get closer than the naive
329    sample mean (again, in terms of absolute error). The absolute errors are provided in supplementary
330    material 3.

331    In terms of confidence/credible interval coverage, the estimators were generally very poor. With the
332    exception of the subsample means, none covered the population mean in either period. The fact that
333    the confidence intervals from the subsamples did cover the population means is not surprising: the
334    subsamples are small ($n = 500$), so the confidence intervals are wide. Of course, increasing the size
335    of the subsamples reduces the width of the confidence intervals, as we show in supplementary
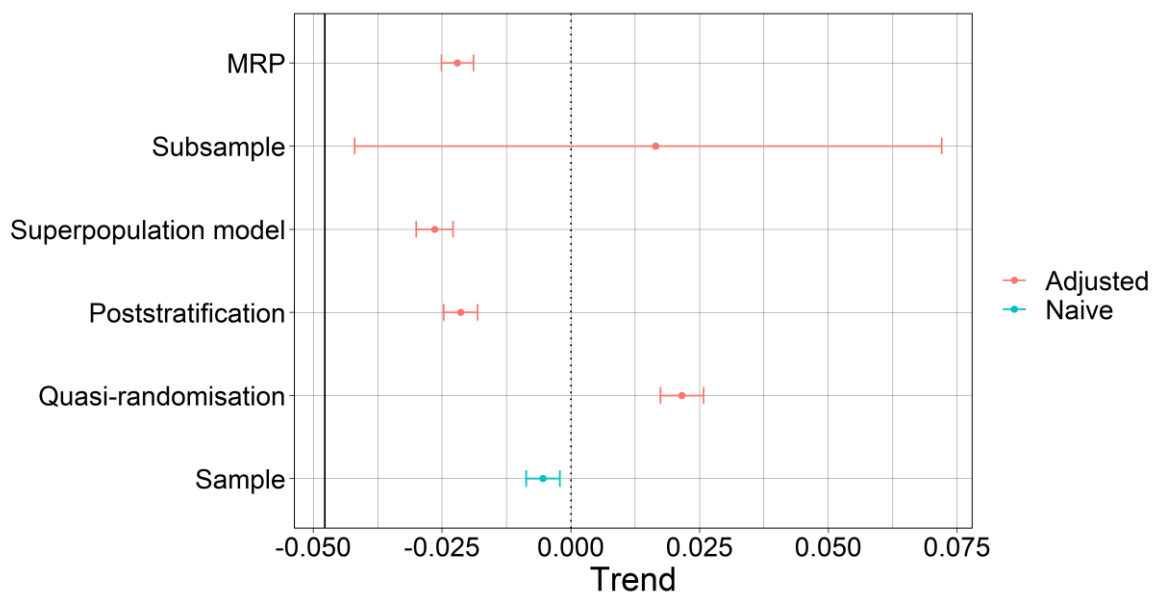336    material 1.

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

## Estimated trend in mean occupancy

Three of the five adjusted point estimates of the trend in mean occupancy are closer than the difference in naive sample means to the true population trends. The other two, the trends from quasi-randomisation and subsampling, are poor. Their point estimates even have the wrong sign. No estimator's credible/confidence interval covers the true trend. The fact that the naïve sample trend underestimates the true trend is a consequence of the time varying representativeness (Bowler et al., 2022; Oliver L. Pescott et al., 2019).
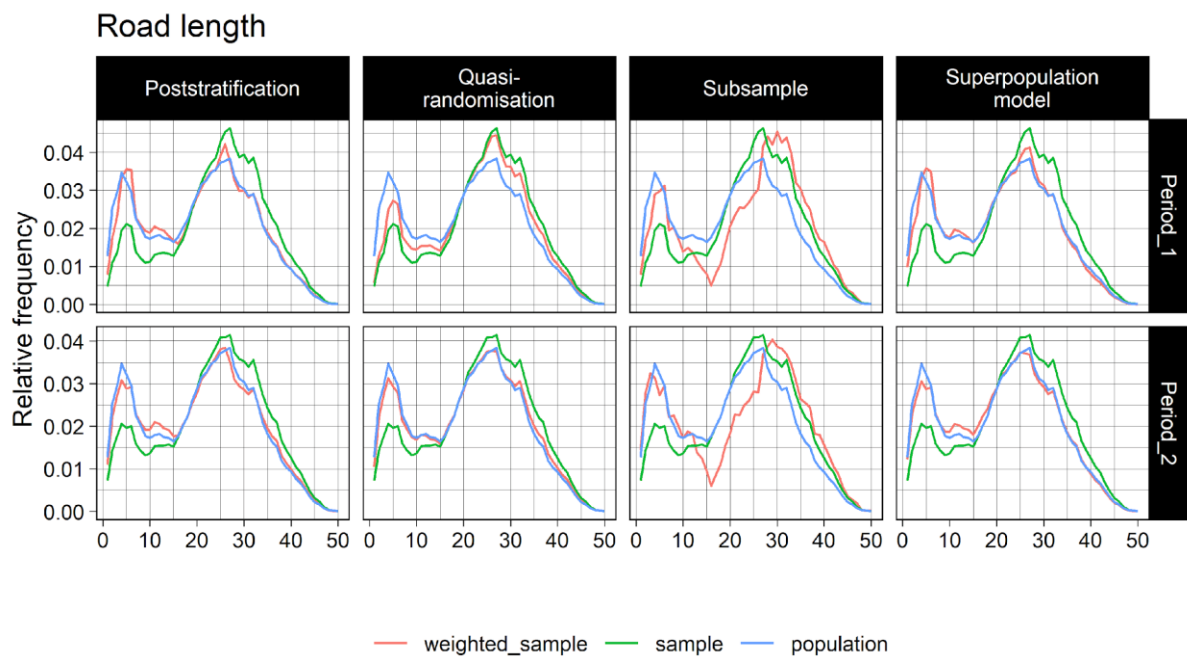
350     Figure 4. Trends in mean occupancy between periods one and two produced by the estimator from
351     each adjustment method, in addition to the naive sample estimate. Error bars delimit 95%
352     confidence/credible intervals. The solid vertical black line denotes the true population trend (-0.047).
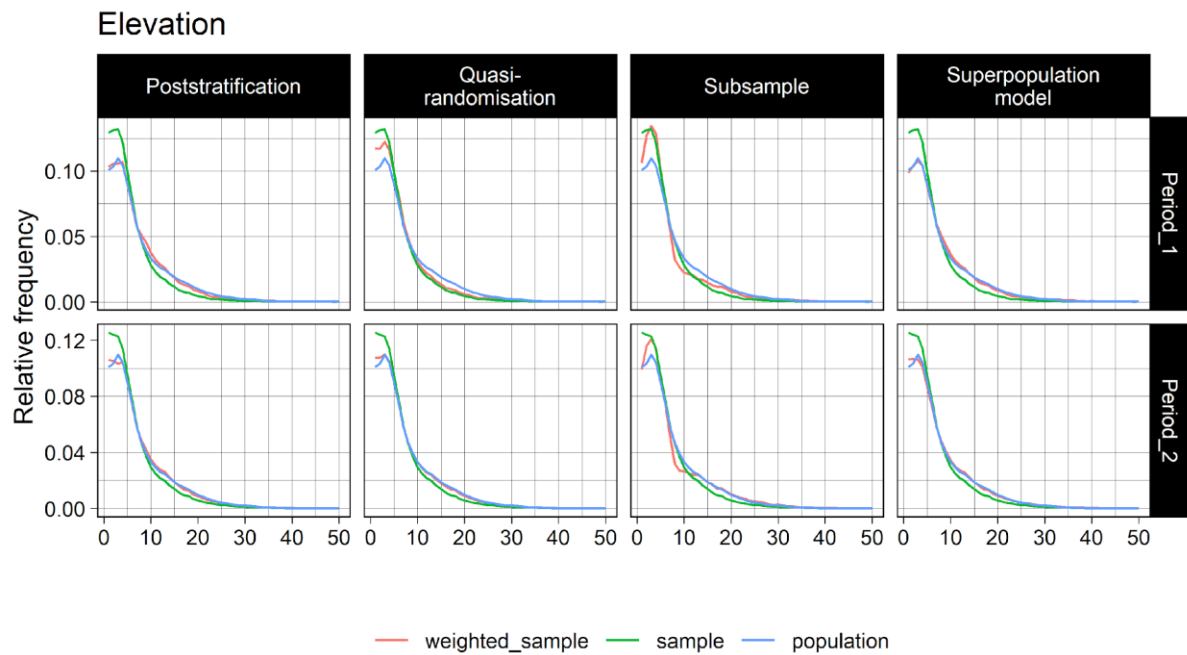
## Distributions of auxiliary variables

354     As measured using Mean Absolute Errors (MAEs), the adjustment methods were generally very good
355     at bringing the distributions of the auxiliaries in the samples closer to those in the population. Figs 5
356     and 6 show the sample and population distributions of two auxiliary variables, road length and
357     elevation, but the MAEs for these and the others can be found in supplementary material 3.
358     Superpopulation modelling and poststratification performed particularly well. Quasi-randomisation
359     offered only a minor improvement in period one. Subsampling was the only approach that did not
360     bring the distributions of the auxiliaries in the sample closer to those in the population.



361

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

Figure 6. Sample, population and weighted sample distributions of the auxiliary variable elevation (Table 1) in periods one and two.

# Discussion

Our experience is that analysts using large, nonprobability samples to monitor biodiversity tend not to account for issues of representativeness. Even where such issues are dealt with, there has been little acknowledgement of the broader panoply of relevant survey sampling methods available to the analyst, no exploration of how these are conceptually (or mathematically) linked and no comparison of their performance in realistic (i.e. relatively data poor) biodiversity monitoring situations. Evidence that a method can work in some discipline, or in simulation studies, is not proof that it will work in all situations. We have demonstrated how such adjustments might be applied using a realistic example of distribution change in a vascular plant over a period of 32 years. This example is realistic in that we do not have access to perfect predictors of occupancy or of sample inclusion. However, it is still likely to be closer to a best-case scenario than otherwise, due to the intense survey effort expended on vascular plants over the British landscape in the recent past (Stroh et al., 2023) and the fact that auxiliary data are relatively accessible in this area.

Our key finding is that the ability to bring the distributions of auxiliary variables in the sample closer to those in the population does not automatically mean that an adjustment will produce a more accurate estimate of a population quantity. For example, poststratification and superpopulation modelling were highly successful at redistributing the auxiliary variables in the samples (Figs 5 and 6). However, this did not translate into large improvements in the estimates of mean occupancy in each time-period or the trend (Figs 3 and 4). It must be the case that the auxiliary variables were not sufficient to describe the key differences between sample and population.

So, what makes a good auxiliary variable? Caughey et al. (2020) listed three criteria: 1) it should predict the response; 2) it should predict sample inclusion; and 3) its distribution in the population should be known. Four of the five auxiliary variables in our example were chosen on the basis that they predict sample inclusion, whereas only one was thought to be a reasonable predictor of the response (occupancy). Whilst it might seem like we prioritised the second criterion over the first, note that we originally included additional predictors of the response. These included soil pH, a known predictor of *C. vulgaris* occupancy (Stroh et al., 2023), and the first and second principal components

394  of climate space in Britain, which we previously found to be important predictors of *C. vulgaris*
395  habitat suitability (Boyd, Harvey, et al., 2023). Including these variables did not improve the estimates
396  of mean occupancy or the trend, as we show in supplementary material 2. We suspect that these
397  variables were redundant, because they are highly correlated with those in Table 1, so it is of little
398  consequence which of these auxiliaries we included.

399  Identifying auxiliary variables that satisfy Caughey and colleagues' (2020) criteria is likely to be the
400  most challenging part of adjusting samples in biodiversity monitoring. In many situations, predictors
401  of the variable of interest and sample inclusion are unknown. Where they are known, data might not
402  be available at the required scale (i.e. their distribution in the population is not known). To illustrate
403  this point, consider the hoverfly *Criorhina asilica*, whose larvae require decaying timber from
404  particular tree species (Stubbs & Falk, 2002). Without data on the locations of those decaying trees, it
405  would likely be impossible to adjust for what is presumably a major determinant of its distribution.
406  For taxa whose habitat requirements are well understood and reflected in available data (e.g. birds),
407  selecting auxiliary variables should be simpler. Nevertheless, in practice, the analyst does not know
408  the truth, so there will always be some guesswork (if this were not the case, statistical modelling
409  would not be required). Transparency regarding availability and choice of auxiliary variables should
410  be an important component of reporting for all biodiversity monitoring.

411  Whilst we are confident that the appropriateness of the auxiliary variables was the limiting factor in
412  our example, it is possible that improvements to the estimators themselves could have improved
413  matters. Potential improvements to MRP are most obvious. For example, interactions between the
414  auxiliary variables could be included in the multilevel model (Ghitza & Gelman, 2013), and multiple
415  time-periods could be modelled at once (Gelman et al., 2018). The question is whether fine-tuning
416  models, which might make them more expensive to run (including interactions in MRP certainly
417  does), is worth marginal gains in accuracy. As Mercer et al., (2018), writing in the context of
418  adjusting survey samples, put it, "[t]he right variables make a big difference for accuracy. Complex
419  statistical methods, not so much."

420  Some have questioned whether it is worth weighting nonprobability samples at all. In opinion polling,
421  for example, there are many examples where weighting or other adjustments did little or nothing to
422  improve the accuracy of inference from nonprobability samples, or even made matters worse (Bailey,
423  2023). In terms of the accuracy of the sample-based estimates, our results suggest that the situation in
424  biodiversity monitoring is similar. Importantly, however, we have also showed that most adjustment
425  methods do what they are supposed to: they turn an unrepresentative sample into a representative one,
426  albeit strictly with respect to the chosen auxiliaries. First principles dictate that, if the auxiliaries are
427  appropriate, this would translate into a more representative sample in terms of the variable of interest
428  and improve the accuracy of inference. We see taxon experts as having a crucial role in identifying
429  appropriate auxiliaries (e.g. Boyd, Harvey, et al., 2023; Smart et al., 2019).

430  It is worth commenting on how we measured the accuracy of the estimated trends. We compared the
431  magnitudes of the estimated and true trends and assessed whether the estimates' confidence/credible
432  intervals covered the true value. Others defined accuracy as the power to "detect" a trend, whereby a
433  method is considered successful if it gets the sign of the trend correct and its uncertainty interval does
434  not span zero (e.g. Valdez et al., 2023). In this power setting, four of the estimators that we considered
435  were able to detect the true trend, including the difference between the naive sample means.

436  We prefer to use the magnitude of the trend as a measure of accuracy for biodiversity monitoring,
437  because many applications in this area are descriptive-inferential, not decision-theoretic (Greenland,
438  2022; Hurlbert et al., 2019; Oliver L. Pescott et al., 2019). That is to say, the final objective of
439  exercises in biodiversity trend estimation is frequently a descriptive indicator, not a binary
440  accept/reject conclusion of change or no change (e.g. Dennis et al., 2019; Powney et al., 2019). The
441  link between Neyman-Pearsonian power and such exercises is often unclear (Amrhein et al., 2019):

442 they are essentially descriptive exercises and as such should be evaluated in terms of the closeness of
443 the sample-based estimate to the truth, not merely in terms of rejecting (typically unrealistic) null
444 hypotheses. The ability to report and consider uncertainty in the trend estimation is essential in
445 making judgements about the risk of bias in biodiversity data (Boyd et al., 2022).

446 Also worth remembering is that we have only applied the adjustments to one species and using a
447 relatively "good" dataset. It is plausible that the adjustments would improve estimates from a smaller
448 or less representative dataset to a greater extent. They will certainly work better for species whose
449 auxiliary variables are easier to identify and reflected in available data.

450 Repeating our analysis with other taxa and datasets would provide a better understanding of in what
451 circumstances we can expect adjustments to perform well. The difficulty will be finding species
452 whose true occupancy (or other variable of interest) is known. One option is to use simulations, but it
453 is extremely important that they are not designed in such a way that the auxiliary variables explain
454 sample inclusion and the variable of interest completely. In this case, the methods will all work very
455 well, but that is not a true reflection of reality.

456 Our concluding message is that statistical adjustments might improve descriptive statistical inference
457 in ecology, but only when combined with expert knowledge and appropriate data. Where there is
458 doubt about the suitability of available auxiliary variables, a safer strategy is to assess the risk of bias
459 qualitatively (Boyd et al., 2022; Meineke & Daru, 2021). If there is deemed to be a risk, it should be
460 reflected in the way that findings are reported (Boyd, Powney, et al., 2023; O L Pescott et al., 2022).
461 This might include using more conservative language and acknowledging that traditional uncertainty
462 intervals are not guaranteed (or even likely) to cover the truth (X.-L. Meng, 2022).

# Acknowledgement

# References

470 Amrhein, V., Trafimow, D., & Greenland, S. (2019). Inferential Statistics as Descriptive Statistics:
471     There Is No Replication Crisis if We Don't Expect Replication. *American Statistician*, *73*(sup1),
472     262–270. https://doi.org/10.1080/00031305.2018.1543137

473 Bailey, M. (2023). *Polling at a crossroads: Rethinking modern survey research*. Cambridge
474     University Press (forthcoming). https://doi.org/10.2753/RSL1061-1975230163

475 Bethlehem, J., Cobben, F., & Schouten, B. (2008). Indicators for the Representativeness of Survey
476     Response. *Statistics Canada's International Symposium Series: Proceedings*, *11*.

477 Bowler, D. E., Callaghan, C. T., Bhandari, N., Henle, K., Barth, M. B., Koppitz, C., Klenke, R.,
478     Winter, M., Jansen, F., Bruelheide, H., & Bonn, A. (2022). Temporal trends in the spatial bias of
479     species occurrence records. *Ecography*. https://doi.org/10.1111/ecog.06219

480 Bowler, D. E., Klaus-, D. E., Conze, J., Suhling, F., Baumann, K., Benken, T., Bönsel, A., Bittner, T.,
481     Drews, A., Günther, A., Isaac, N. J. B., & Petzold, F. (2021). Winners and losers over 35 years
482     of dragonfly and damselfly distributional change in Germany. *Diversity and Distributions*,
483     *27*(August 2020), 1353–1366. https://doi.org/10.1111/ddi.13274

484 Boyd, R., Harvey, M., Roy, D., Barber, T., Haysom, K., & ... (2023). Causal inference and large-scale
485     expert validation shed light on the drivers of SDM accuracy and variance. *Diversity and*

*Distributions*, 1–11. https://doi.org/10.1111/ddi.13698

Boyd, R., Powney, G. D., Burns, F., Danet, A., Duchenne, F., Grainger, M. J., Jarvis, S. G., Martin, G., Nilsen, E. B., Porcher, E., Stewart, G. B., Wilson, O. J., & Pescott, O. L. (2022). ROBITT: A tool for assessing the risk-of-bias in studies of temporal trends in ecology. *Methods in Ecology and Evolution*, *13*(March), 1497– 1507. https://doi.org/10.1111/2041-210X.13857

Boyd, R., Powney, G. D., & Pescott, O. L. (2023). We need to talk about nonprobability samples. *Trends in Ecology & Evolution*, *xx*(xx), 1–11. https://doi.org/10.1016/j.tree.2023.01.001

Caughey, D., Berinsky, A., Chatfield, S., Hartman, E., Schickler, E., & Sekhon, J. (2020). *Target Estimation and Calibration Weighting for Unrepresentative Survey Samples*. Cambridge University Press.

Dennis, E. B., Brereton, T. M., Morgan, B. J. T., Fox, R., Shortall, C. R., Prescott, T., & Foster, S. (2019). Trends and indicators for quantifying moth abundance and occupancy in Scotland. *Journal of Insect Conservation*, *23*(2), 369–380. https://doi.org/10.1007/s10841-019-00135-z

Fick, S. E., & Hijmans, R. J. (2017). WorldClim 2 : new 1-km spatial resolution climate surfaces for global land areas. *International Journal of Climatology*. https://doi.org/10.1002/joc.5086

Geldmann, J., Heilmann-Clausen, J., Holm, T. E., Levinsky, I., Markussen, B., Olsen, K., Rahbek, C., & Tøttrup, A. P. (2016). What determines spatial bias in citizen science? Exploring four recording schemes with different proficiency requirements. *Diversity and Distributions*, *22*(11), 1139–1149. https://doi.org/10.1111/ddi.12477

Gelman, A. (2007). Struggles with survey weighting and regression modeling. *Statistical Science*, *22*(2), 153–164. https://doi.org/10.1214/088342306000000691

Gelman, A., & Little, T. (1997). poststratification into many categories using hierarchical regression. *Survey Methodology*, *23*(2), 127–335.

Gelman, A., Phillips, J., Lax, J., Gabry, J., & Trangucci, R. (2018). *Using Multilevel Regression and Poststratification to Estimate Dynamic Public Opinion*. http://www.stat.columbia.edu/~gelman/research/unpublished/MRT(1).pdf

Ghitza, Y., & Gelman, A. (2013). Deep interactions with MRP: Election turnout and voting patterns among small electoral subgroups. *American Journal of Political Science*, *57*(3), 762–776. https://doi.org/10.1111/ajps.12004

Greenland, S. (2022). Divergence versus decision P-values: A distinction worth making in theory and keeping in practice: Or, how divergence P-values measure evidence even when decision P-values do not. *Scandinavian Journal of Statistics*, *50*(1), 54–88. https://doi.org/10.1111/sjos.12625

Gregory, R. D., Van Strien, A., Vorisek, P., Meyling, A. W. G., Noble, D. G., Foppen, R. P. B., & Gibbons, D. W. (2005). Developing indicators for European birds. *Philosophical Transactions of the Royal Society B: Biological Sciences*, *360*(1454), 269–288. https://doi.org/10.1098/rstb.2004.1602

Hudson, L. N., Newbold, T., Contu, S., Hill, S. L. L., Lysenko, I., Palma, D., Phillips, H. R. P., Senior, R. A., Bennett, D. J., Booth, H., Garon, M., Michelle, L., Correia, D. L. P., Day, J., Echeverr, S., Harrison, K., Ingram, D. J., Jung, M., Kemp, V., … Fernando, A. B. (2014). *The PREDICTS database : a global database of how local terrestrial biodiversity responds to human impacts*. 4701–4735.

Hurlbert, S. H., Levine, R. A., & Utts, J. (2019). Coup de Grâce for a Tough Old Bull: "Statistically Significant" Expires. *American Statistician*, *73*(sup1), 352–357. https://doi.org/10.1080/00031305.2018.1543616

Intermap. (2009). *NEXTMap British Digital Terrain 50m resolution (DTM10) Model Data by Intermap*. NERC Earth Observation Centre. https://catalogue.ceda.ac.uk/uuid/f5d41db1170f41819497d15dd8052ad2

Johnston, A., Matechou, E., & Dennis, E. B. (2022). Outstanding challenges and future directions for biodiversity monitoring using citizen science data. *Methods in Ecology and Evolution*, *February*. https://doi.org/10.1111/2041-210X.13834

Johnston, A., Moran, N., Musgrove, A., Fink, D., & Baillie, S. R. (2020). Estimating species distributions from spatially biased citizen science data. *Ecological Modelling*, *422*(December 2019), 108927. https://doi.org/10.1016/j.ecolmodel.2019.108927

Lohr, S. (2022). *Sampling: Design and analysis* (3rd ed.). CRC Press.

Lumley, T. (2010). *Complex surveys: A Guide to Analysis Using R* (1st ed.). Wiley.

Makela, S., Si, Y., & Gelman, A. (2014). Statistical Graphics for Survey Weights. *Revista Colombiana de Estadística*, *37*(2Spe), 285–295. https://doi.org/10.15446/rce.v37n2spe.47937

Meineke, E. K., & Daru, B. H. (2021). Bias assessments to expand research harnessing biological collections. *Trends in Ecology and Evolution*, *36*(12), 1071–1082. https://doi.org/10.1016/j.tree.2021.08.003

Meng, X.-L. (2022). Double Your Variance, Dirtify Your Bayes, Devour Your Pufferfish, and Draw your Kidstrogram. *The New England Journal of Statistics in Data Science*, *0*, 1–20. https://doi.org/10.51387/22-nejsds6

Meng, X. (2022). Comments on the Wu ( 2022 ) paper by Xiao-Li Meng 1 : Miniaturizing data defect correlation : A versatile strategy for handling non-probability samples. *Survey Methodology*, *48*(2), 1–22.

Meng, X. L. (2018). Statistical paradises and paradoxes in big data (I): Law of large populations, big data paradox, and the 2016 us presidential election. *Annals of Applied Statistics*, *12*(2), 685–726. https://doi.org/10.1214/18-AOAS1161SF

Mercer, A., Lau, A., & Kennedy, C. (2018). For Weighting Online Opt-In Samples, What Matters Most? In *Pew Research Center* (pp. 1–55). https://pewrsr.ch/3heqknn

Morton, R., Marston, C., O'Neil, A., & Rowland, C. (2022). *Land Cover Map 2018 (1km summary rasters, GB and N. Ireland)*. NERC EDS Environmental Information Data Centre. https://doi.org/https://doi.org/10.5285/9b68ee52-8a95-41eb-8ef1-8d29e2570b00

Nelson, G., & Ellis, S. (2019). The history and impact of digitization and digital data mobilization on biodiversity research. *Philosophical Transactions of the Royal Society B: Biological Sciences*, *374*(1763), 2–10. https://doi.org/10.1098/rstb.2017.0391

ONS. (2021). *UK Office for National Statistics Postcode Directory*. https://geoportal.statistics.gov.uk/datasets/ons-postcode-directory-november-2022/about

Outhwaite, C., Gregory, R. D., Chandler, R. E., Collen, B., & Isaac, N. J. B. (2020). Complex long-term biodiversity change among invertebrates, bryophytes and lichens. *Nature Ecology & Evolution*. https://doi.org/10.1038/s41559-020-1111-z

Perring, F., & Walters, S. (1962). *Atlas of the British flora*. Thomas Nelson & sons.

Pescott, O L, Stroh, P. A., Humphrey, T. A., & Walker, K. J. (2022). Simple methods for improving the communication of uncertainty in species ' temporal trends. *Ecological Indicators*, *141*(May). https://doi.org/https://doi.org/10.1016/j.ecolind.2022.109117

Pescott, Oliver L., Humphrey, T. A., Stroh, P. A., & Walker, K. J. (2019). Temporal changes in distributions and the species atlas: How can British and Irish plant data shoulder the inferential

575       burden? *British & Irish Botany*, *1*(4), 250–282. https://doi.org/10.33928/bib.2019.01.250

576  Powney, G. D., Carvell, C., Edwards, M., Morris, R. K. A., Roy, H. E., Woodcock, B. A., & Isaac, N.
577       J. B. (2019). Widespread losses of pollinating insects in Britain. *Nature Communications*,
578       *10*(2019), 1–6. https://doi.org/10.1038/s41467-019-08974-9

579  Preston, C.D., Pearman, D.A. & Dines, T. D. (2002). *New Atlas of the British and Irish Flora.* (eds).
580       Oxford University Press.

581  Rowland, C., Marston, C., Morton, R., & O'Neil, A. (2020). *Land Cover Map 1990 (1km dominant*
582       *target class, GB) v2.* NERC EDS Environmental Information Data Centre.
583       https://doi.org/https://doi.org/10.5285/f5e3bd00-efd0-4dc6-a454-aa597d84764a

584  Smart, S. M., Jarvis, S. G., Mizunuma, T., Herrero-Jáuregui, C., Fang, Z., Butler, A., Alison, J.,
585       Wilson, M., & Marrs, R. H. (2019). Assessment of a large number of empirical plant species
586       niche models by elicitation of knowledge from two national experts. *Ecology and Evolution*,
587       *9*(22), 12858–12868. https://doi.org/10.1002/ece3.5766

588  Stroh, P. A., Walker, K., Humphrey, T. A., Pescott, O. L., & Burkmar, R. (2023). *Plant Atlas 2020:*
589       *Mapping Changes in the Distribution of the British and Irish Flora*. Princeton Univ. Press.

590  Stubbs, A., & Falk, S. (2002). *British Hoverflies*. British Entomological and Natural History Society.

591  UNEP-WCMC, & IUCN. (2020). *Protected Planet: The World Database on Protected Areas*
592       *(WDPA)/The Global Database on Protected Areas Management Effectiveness*.
593       https://www.protectedplanet.net/en/thematic-areas/wdpa

594  Valdez, J. W., Callaghan, C. T., Junker, J., Purvis, A., Hill, S. L. L., & Pereira, H. M. (2023). The
595       undetectability of global biodiversity trends using local species richness. *Ecography*, *2023*(3), 1–
596       14. https://doi.org/10.1111/ecog.06604

597  Valliant, R. (2020). Comparing Alternatives for Estimation from Nonprobability Samples. *Journal of*
598       *Survey Statistics and Methodology*, *8*(2), 231–263. https://doi.org/10.1093/jssam/smz003

599  Valliant, R., Dever, J. A., & Kreuter, F. (2018). *Practical tools for designing and weighting survey*
600       *samples* (2nd ed.). Springer Cham. https://doi.org/https://doi.org/10.1007/978-3-319-93632-1

601  van Strien, A. J., & van Grunsven, R. H. A. (2023). In the past 100 years dragonflies declined and
602       recovered by habitat restoration and climate change. *Biological Conservation*, *277*(December
603       2022), 109865. https://doi.org/10.1016/j.biocon.2022.109865

604  Van Swaay, C.A.M., Plate, C. L., & Van Strien, A. J. (2002). Monitoring butterflies in the
605       Netherlands: how to get unbiased indices. *Proceedings of the Section Experimental and Applied*
606       *Entomology of the Netherlands Entomological Society*, *13*, 21–27.

607  Van Swaay, Chris A.M., Nowicki, P., Settele, J., & Van Strien, A. J. (2008). Butterfly monitoring in
608       Europe: Methods, applications and perspectives. *Biodiversity and Conservation*, *17*(14), 3455–
609       3469. https://doi.org/10.1007/s10531-008-9491-4

610  Weiser, E. L., Diffendorfer, J. E., Lopez-Hoffman, L., Semmens, D., & Thogmartin, W. E. (2020).
611       Challenges for leveraging citizen science to support statistically robust monitoring programs.
612       *Biological Conservation*, *242*(October 2019). https://doi.org/10.1016/j.biocon.2020.108411

613  Wu, C., & Sitter, R. R. (2001). A model-calibration approach to using complete auxiliary information
614       from survey data. *Journal of the American Statistical Association*, *96*(453), 185–193.
615       https://doi.org/10.1198/016214501750333054

616