

On the trade-off between accuracy and spatial resolution when estimating species occupancy from geographically biased samples

*Robin J. Boyd, Diana E. Bowler, Nick J. B. Isaac and Oliver L. Pescott

UK Centre for Ecology and Hydrology, Benson Ln., Wallingford, Oxfordshire, UK, OX10 8BB

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

Orcid IDs: OLP- 0000-0002-0685-8046; NJBI 0000-0002-4869-8052

Abstract

Species occupancy is often defined as the proportion of areal units (sites) in a landscape that the focal species occupies, but it is usually estimated from the subset of sites that have been sampled. Assuming no measurement error, we show that three quantities—the degree of sampling bias (in terms of site selection), the proportion of sites that have been sampled and the variability of true occupancy across sites—determine the extent to which a sample-based estimate of occupancy differs from its true value across the wider landscape. That these are the only three quantities (measurement error notwithstanding) to affect the accuracy of estimates of species occupancy is the fundamental insight of the “Meng equation”, an algebraic re-expression of statistical error. We use simulations to show how each of the three quantities vary with the spatial resolution of the analysis and that absolute estimation error is lower at coarser resolutions. Absolute error scales similarly with resolution regardless of the size and clustering of the virtual species’ distribution. Finely resolved estimates of species occupancy have the potential to be more useful than coarse ones, but this potential is only realised if the estimates are at least reasonably accurate. Consequently, wherever there is the potential for sampling bias, there is a trade-off between spatial resolution and accuracy, and the Meng equation provides a theoretical framework in which analysts can consider the balance between the two. An obvious next step is to consider the implications of the Meng equation for estimating a time trend in species occupancy, where it is the confounding of error and true change that is of most interest.

Key words: sampling bias; spatial grain; representativeness; biodiversity monitoring; spatial pattern

1 Introduction

The proportion of “sites” occupied by some species (its occupancy) is often of interest to ecologists (Kéry & Royle, 2016). Sites were originally conceived as discrete habitat patches or relatively small sampling units (MacKenzie et al., 2002), but increasingly they represent contiguous larger-scale units defined by the analyst (e.g. squares on a map; Van Strien et al., 2013). This latter definition has often been used when estimating species occupancy at national and supranational scales (Boyd, August, et al., 2023; Isaac et al., 2014).

In most circumstances—and particularly at fine scales across large areas—data are not available for all sites, so occupancy must be estimated from the subset of sites that have been sampled (Kéry & Royle, 2016). If the focal species is more or less likely to occupy sampled than non-sampled sites, then the sample is geographically biased (a formal definition is provided below), and the sample-based estimate will differ from its true value across the wider landscape (Boyd, Powney, et al., 2023; Meng, 2018). Geographic sampling biases are just one source of error when estimating species occupancy, the other major source being measurement error at sampled sites (MacKenzie et al., 2002).

42 A further complication when estimating species occupancy is that it varies with spatial resolution.
43 Occupancy always increases as the resolution is coarsened, but the rate at which it increases depends
44 on the size and clustering of the species' distribution at the finer scales (Azaele et al., 2012; Kunin,
45 1998; Wilson et al., 2004). All else being equal, fine scale estimates of species occupancy are
46 preferable to coarse ones. For example, colonisations and local extinctions at small-scale sites are
47 more probable than at larger scales, so working at a finer resolution means that occupancy is more
48 sensitive to change (Dennis et al., 2019).

49 Although estimates of occupancy are nominally more useful at fine scales, there are reasons to work
50 at coarser resolutions too. One reason is that, given finite resources, sampling at a fine scale might
51 come at the expense of sampling over a large geographic area. Another is that the effects of sampling
52 bias become more pronounced where there are more sites in the landscape (Boyd, Powney, et al.,
53 2023; Meng, 2018), which is obviously the case at finer resolutions (i.e. where the sites are smaller).
54 The fact that sampling biases are more pervasive at finer spatial resolutions raises questions about
55 how the accuracy of estimates of species occupancy scales with resolution. Although working at
56 coarser resolutions will clearly improve accuracy at the extremes—we can be surer a species occupies
57 planet Earth than a set of small plots on its surface—how accuracy varies along the gradient from fine
58 to coarse resolutions under sampling bias has not, to our knowledge, been investigated in ecology.

59 Here then, we investigate how the error of sample-based estimators of species occupancy vary with
60 spatial resolution. Assuming no false absences (or that a model has adequately corrected them), we
61 begin by demonstrating that three, and only three, quantities determine the magnitude of the error: the
62 degree of sampling bias (in terms of site selection), the proportion of sites sampled and the variability
63 of true occupancy across sites. That these are the only quantities affecting estimation error is a key
64 implication of Meng's (2018) decomposition of survey error. We use simulations to show how each
65 of the three quantities, and both relative and absolute error, vary with spatial resolution under
66 sampling bias (at the finest resolution) and how varying the level of sampling bias affects the error. A
67 trade-off emerges between finely resolved and accurate estimates, which we discuss in detail.

68 2 Methods

69 2.1 Quantifying estimation error

70 We consider a landscape comprising N contiguous sites of equal area. The presence of at least one
71 individual of the focal species is a binary variable Y taking the value 1 at sites where it is present and
72 0 elsewhere. Occupancy $P(Y = 1)$ is the proportion of sites at which the species is present, which is
73 equivalent to the mean of Y across sites \bar{Y} . Of the N sites, a subset n are sampled. Whether each site is
74 one of the n sampled sites is another binary variable R ($R = 1$ where the site is sampled and $R = 0$
75 otherwise). It is not possible to calculate mean occupancy across all N sites, \bar{Y}_N , because information
76 on Y is not available for sites with $R = 0$. Instead, it is common to *estimate* \bar{Y}_N as mean occupancy
77 across sampled sites \bar{Y}_n .

78 Assuming no measurement errors, or that a model has corrected them, the absolute error of \bar{Y}_n as an
79 estimator of \bar{Y}_N is (Meng, 2018)

$$\bar{Y}_n - \bar{Y}_N = \rho(R, Y) \sqrt{\frac{1-f}{f}} \sigma_Y. \quad \text{equation 1}$$

80 The first quantity on the right, $\rho(R, Y)$, is the (population) correlation between Y and R . It is a
81 measure of both the sign and magnitude of *sampling bias*. In simple terms, $\rho(R, Y)$ is negative where
82 Y is generally smaller in the sample than in the population and vice versa. f is the sampling rate
83 (n/N), and the second quantity on the right is a measure of *data quantity*. The final quantity σ_Y is the
84 population standard deviation of Y . It is 0 where Y is constant, in which case a sample size of 1 is

85 sufficient to estimate \bar{Y}_N with no error, and it is largest where Y is most variable. Hence, it can be
86 considered a measure of “*problem difficulty*” (Meng, 2018), although we refer to it as occupancy
87 variability given the context in which we are working.

88 Importantly, eq. 1 gives the absolute error of \bar{Y}_n as an estimator of \bar{Y}_N for a given sample: that is, for
89 one realisation of R . In what follows, we consider replicate realisations of R from given R -generating
90 (i.e. sampling) mechanisms and the average $\bar{Y}_n - \bar{Y}_N$ across those samples.

91 2.2 Effects of spatial resolution on error

92 Eq. 1 provides a basis for understanding the effects of resolution on absolute error when estimating
93 species occupancy. Assuming perfect detection, it implies that there are three, and only three, ways to
94 reduce error: decrease the sampling bias $\rho(R, Y)$, increase the sampling rate f and/or decrease the
95 occupancy variability σ_Y . Below we describe a set of simulations that demonstrate the effects of
96 coarsening the spatial resolution on each of these quantities and on both absolute and relative error.

97 2.3 Simulation setup

98 2.3.1 Virtual landscape, species and samples

99 The virtual landscape comprises a square grid of $N = 6400$ cells (80×80) at the finest resolution.
100 Each cell might represent, say, a 1×1 km grid square, but the precise definition is not important for
101 drawing general conclusions.

102 We simulated six species’ geographic distributions of different sizes and with different levels of
103 clustering in the virtual landscape. Our approach was a simplified version of the one used by (Guélat
104 & Kéry, 2018). For each species, the first step was to populate every cell in the landscape with a
105 continuous index X sampled from a multivariate normal distribution

$$X \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\delta}), \quad \text{equation 2}$$

106 where $\boldsymbol{\mu}$ is an N -vector of zeros (i.e. mean X for each grid cell) and $\boldsymbol{\delta}$ is an $N \times N$ covariance matrix.
107 We used an exponential decay function to define the covariance matrix

$$\boldsymbol{\delta} = e^{-\varphi \mathbf{D}_{ij}}, \quad \text{equation 3}$$

108 where φ is the decay constant and \mathbf{D}_{ij} is the Euclidian distance between grid cells i and j . Larger
109 values of φ result in patchier distributions, because the covariance between grid cells diminishes
110 faster with the distance between them.

111 The next step was to convert the continuous index X to a binary one (i.e. occupied vs unoccupied)
112 with a specified proportion of cells being occupied. For each species, we set a threshold percentile of
113 X across grid cells ($1 - \bar{Y}_N$) above which the cell was designated occupied and below which it was
114 designated unoccupied. Table 1 lists the parameters used to simulate each species’ geographic
115 distribution and the resulting properties of those distributions.

116 It was important that the simulated species’ distributions spanned a range of plausible sizes and levels
117 of clustering, because these properties determine how \bar{Y}_N scales with resolution (Kunin, 1998). We
118 tested whether the distributions covered sufficiently wide ranges of these parameters using their
119 fractal dimensions (Kunin, 1998). The fractal dimension D of a species’ distribution is given by $D =$
120 $2(1 - b)$, where b is the slope of its scale-area curve or occupancy-area relationship (Hartley &
121 Kunin, 2003). We calculated b over the finest three resolutions, because, for the medium and common
122 species, including the coarsest two resolutions resulted in nonlinear scale-area curves (i.e. their
123 distributions are non-fractal at coarse scales). The theoretical limits of the fractal dimension are 0,
124 representing a species whose distribution is very sparse, and 2, representing a species whose
125 distribution is very clustered (Hartley & Kunin, 2003). Our virtual species’ distributions spanned most

126 of this range (0.31–1.64). Like (Wilson et al., 2004), we found that D is positively related to \bar{Y}_N ,
 127 which reflects the facts that a small distribution can only be so clustered, and a large distribution can
 128 only be so dispersed.

129 Table 1. Properties of the six virtual species’ distributions at the finest spatial resolution. The
 130 autocorrelation parameter is the exponential decay constant in eq. 3, and higher values produce a more
 131 dispersed distribution. The theoretical limits for the fractal dimension are 0, representing a highly
 132 dispersed species, and 2, representing a very clustered one. The fractal dimension also varies with \bar{Y}_N
 133 (Wilson et al., 2004).

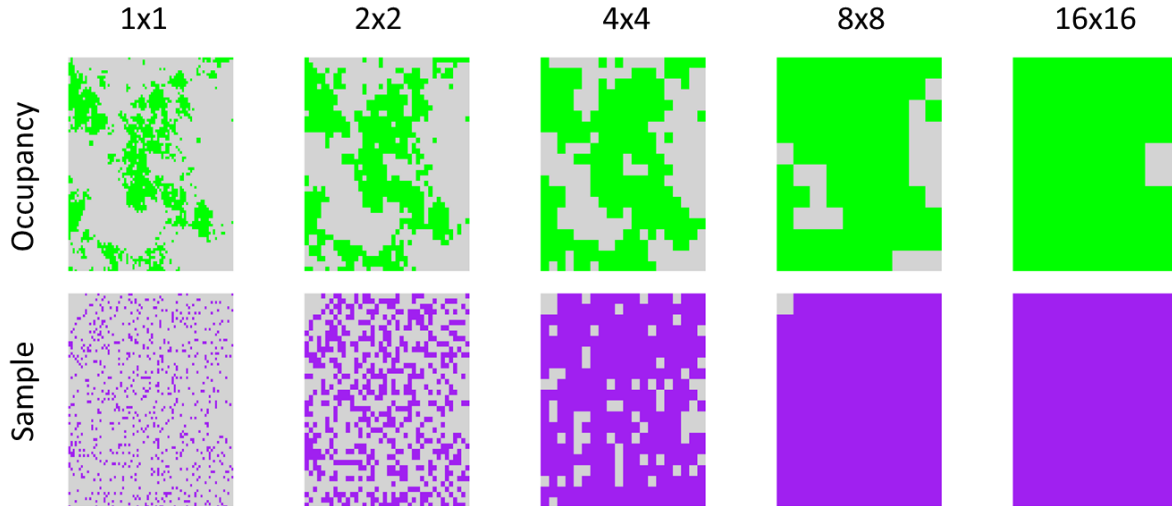
Distribution properties	Exponential decay parameter in autocorrelation function	Proportion of sites occupied (at the finest scale)	Fractal dimension
Rare and sparse	0.6	0.01	0.31
Rare and clustered	0.1	0.01	0.87
Medium and sparse	0.6	0.25	1.20
Medium and clustered	0.1	0.25	1.39
Common and sparse	0.6	0.5	1.57
Common and clustered	0.1	0.5	1.64

134

135 For each species, we simulated 100 virtual samples at the finest resolution. Whilst it might seem more
 136 logical to simulate one set of samples for all species, this would not allow control over $\rho(R, Y)$, the
 137 sampling bias, which depends on the focal species’ geographic distribution. For most simulations, we
 138 simulated the samples in such a way that $E_R[\rho(R, Y)] \sim 0.05$ and $f = 0.1$, where $E_R[\rho(R, Y)]$ is the
 139 expectation (average) of $\rho(R, Y)$ over the 100 simulated samples (i.e. with respect to R). See the
 140 supplementary Fig. S1 for the distributions of $\rho(R, Y)$ across samples for each species. We based the
 141 values of $\rho(R, Y)$ and f on an empirical example: a citizen science dataset on vascular plant sampling
 142 and the species *Calluna vulgaris*’ occupancy in Britain (Boyd et al., 2023). Whilst we generally set
 143 $E_R[\rho(R, Y)] \sim 0.05$ and $f = 0.1$, we also demonstrate the effects of varying both parameters (in the
 144 supplementary material for f). Switching the sign of $\rho(R, Y)$ (i.e. whether occupancy is larger or
 145 smaller in the sample than the population) would switch the sign of the error in the estimate of mean
 146 occupancy, but for simplicity we only present the positive case.

147 2.4 Analysis of error at each resolution

148 The goal of our analysis was to determine how the absolute error of \bar{Y}_n as an estimator of \bar{Y}_N ($\bar{Y}_n -$
 149 \bar{Y}_N ; assuming perfect detection) varies with spatial resolution. Starting at the finest resolution, we
 150 calculated the value of each quantity in eq. 1 (including the absolute error; averaged across the 100
 151 samples). We then coarsened the resolution by aggregating every square of four grid cells into one
 152 (i.e. doubling the length and width of the site). After coarsening the resolution, we recalculated each
 153 quantity in eq. 1, coarsened the resolution again and repeated the process until each grid cell was $16 \times$
 154 its original height and width. Fig. 1 shows how a species’ distribution (medium and clustered; Table
 155 1) and a sample vary with resolution.



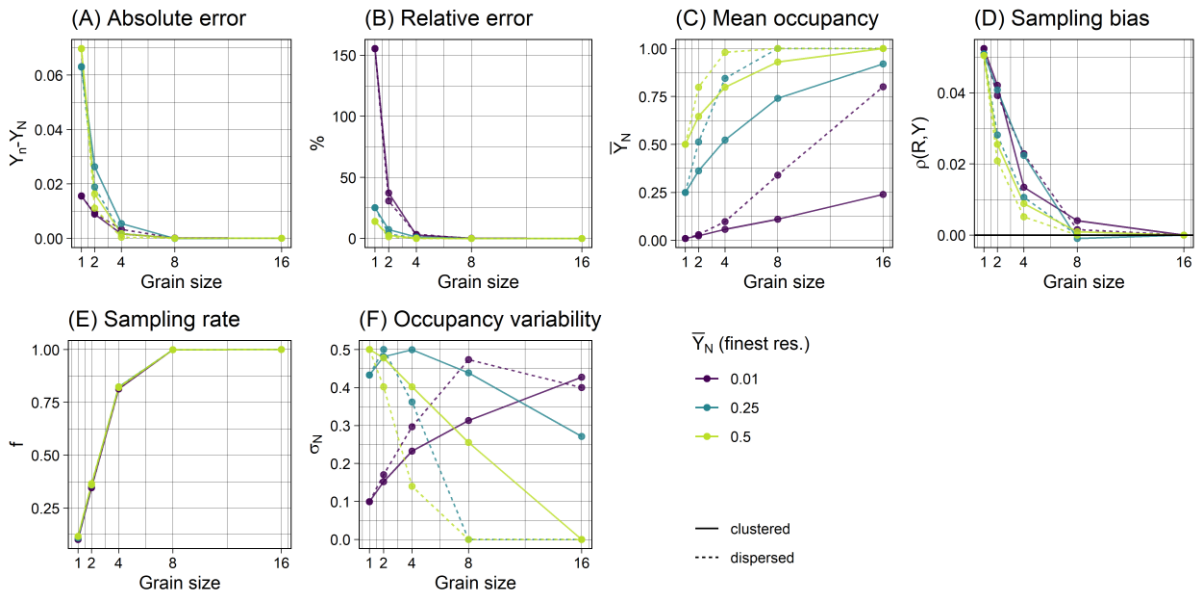
156

157 Figure 1. Top row: a virtual species' ("medium and clustered"; Table 1) geographic distribution at
 158 each spatial resolution. Green cells are occupied, and grey cells are not. Bottom row: a virtual sample
 159 at each resolution. $\rho(R, Y) \sim 0.05$ and $f \sim 0.1$ at the finest resolution (1×1). Purple cells are sampled,
 160 and grey cells are not. Sampled cells may be either occupied or not.

161 3 Results

162 3.1 Error

163 For all virtual species, estimates of occupancy are more accurate at coarser resolutions. This result is
 164 evident both in terms of the absolute actual error (Fig. 2A), which is on the left side of eq. 1, and the
 165 relative actual error (Fig. 2B), which expresses the absolute error as a percentage of true occupancy.
 166 Relative error is larger for rare species. Absolute error is larger for the medium and common species,
 167 particularly at the finer resolutions.



168

169 Figure 2. (A) absolute error, (B) relative error (i.e. the absolute error expressed as a percentage of true
 170 occupancy), (C) mean occupancy (i.e. true occupancy), (D) sampling bias, (E) sampling rate and (F)
 171 occupancy variability σ_Y at each resolution. The resolution is the height and width of the grid cells in
 172 arbitrary units. Points represent the average of each statistic over 100 simulated samples. At the finest
 173 resolution, $\rho(R, Y) \sim 0.05$ and $f \sim 0.1$, the target values for the simulations.

174 3.2 True occupancy

175 Although well-documented (Azaele et al., 2012; Kunin, 1998), it is worth revisiting the scaling
176 properties of \bar{Y}_N (i.e. a species' true occupancy) here, because they provide insight into the scaling
177 properties of error. \bar{Y}_N always increases with resolution, but the rate at which it increases depends on
178 the properties of the species' distribution at the finest resolution (Fig. 2C). Species that are common
179 and sparsely distributed at the finest resolution quickly reach $\bar{Y}_N = 1$ as the resolution is coarsened.
180 By contrast, species that are rare and clustered at the finest resolution do not reach $\bar{Y}_N = 1$ at any of
181 the resolutions we considered (Fig. 2A).

182 3.3 Sampling bias

183 In our simulations, the sampling bias $\rho(R, Y)$ tends towards 0 as the resolution is coarsened. There are
184 plausible scenarios in which it will not, however, a point that we expand on in the Discussion.

185 3.4 Sampling rate

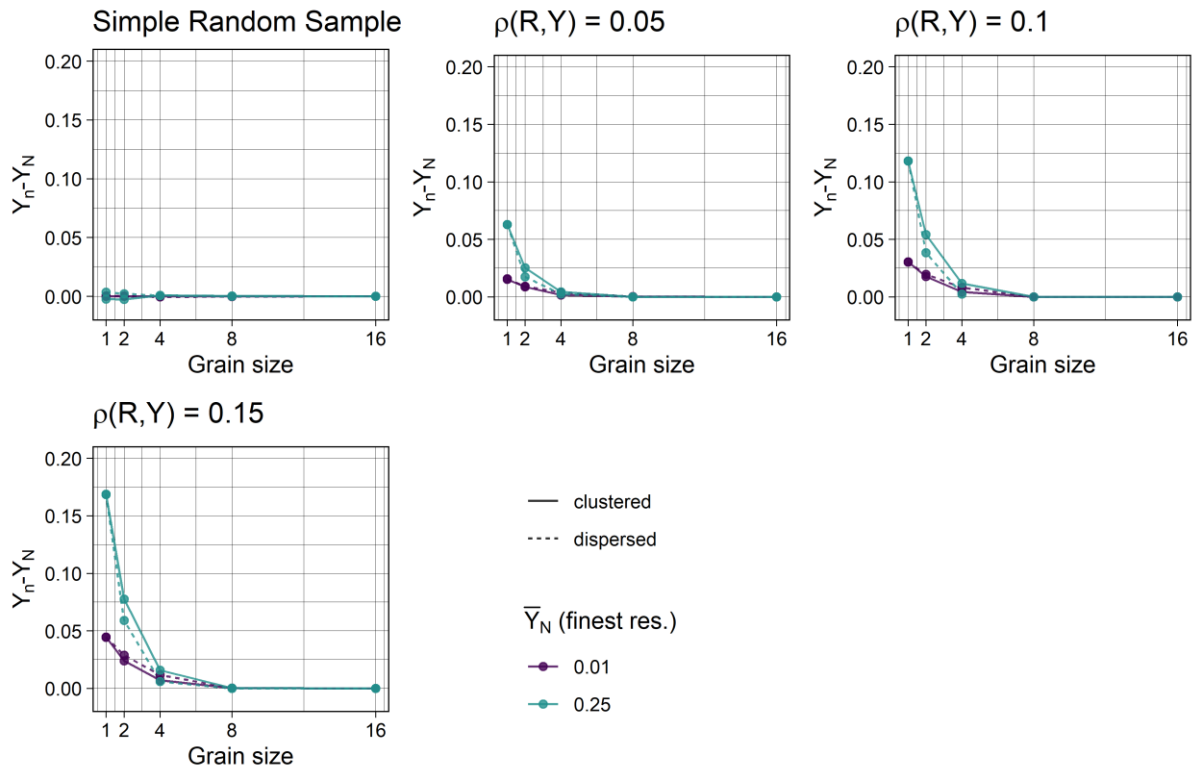
186 The sampling rate f scales in a similar way with resolution to \bar{Y}_N . It always increases with resolution,
187 and sparser samples increase at a greater rate. In our simulations, sparser samples are slightly more
188 likely for the sparsely distributed species because of the correlation between the species' distributions
189 and the samples (i.e. the programmed sampling bias). Hence, f does not increase at exactly the same
190 rate for all species.

191 3.5 Occupancy variability

192 As occupancy is binary, its standard deviation σ_Y is given by $\sqrt{\bar{Y}_N(1 - \bar{Y}_N)}$. σ_Y is largest where \bar{Y}_N is
193 near 0.5 and smallest where \bar{Y}_N is near 0 or 1. Given that \bar{Y}_N increases with resolution (Fig. 2C),
194 coarsening the resolution for species with $\bar{Y}_N < 0.5$ increases σ_Y until $\bar{Y}_N = 0.5$ (Fig. 2F). Further
195 coarsening the resolution decreases σ_Y , because \bar{Y}_N moves away from 0.5 and towards 1. For species
196 with $\bar{Y}_N \geq 0.5$ at the finest resolution, coarsening the resolution always decreases σ_Y .

197 3.6 Scaling of error with resolution at different levels of sampling bias

198 In most simulations, we set $\rho(R, Y) \sim 0.05$ at the finest resolution, but it is instructive to see how
199 actual error scales with resolution under different levels of sampling bias. Error generally scales in the
200 same way with resolution regardless of the level of sampling bias, but is greater in magnitude under
201 stronger sampling bias (Fig. 3). Under a simple random sample at the finest resolution, where the
202 expected sampling bias $E_R[\rho(R, Y)] = 0$, there is roughly no error at any resolution (recalling that we
203 present the average error across samples, which essentially removes sampling error). Note that we
204 were not able to simulate highly biased samples ($E_R[\rho(R, Y)] \sim 0.15$) for the common species (blue
205 lines in Fig. 3). For these species, \bar{Y}_N is very different to f , which makes a large and positive $\rho(R, Y)$
206 highly unlikely, and our algorithm for generating the samples could not achieve it.



207

208 Figure 3. Absolute error at each resolution under four levels of sampling bias $\rho(R, Y)$ (at the finest
 209 resolution). The resolution is the height and width of the grid cells in arbitrary units. The simple
 210 random sample has approximately no sampling bias at the finest resolution. Each line represents one
 211 virtual species: red = rare, green = medium and blue = common. Solid lines depict species with
 212 clustered distributions at the finest resolution and dashed lines indicate species that are highly
 213 dispersed at the finest resolution. Points represent the average of each statistic over 100 simulated
 214 samples. $f \sim 0.1$ at the finest resolution in all cases.

215 4 Discussion

216 Nobody would dispute the fact that estimates of species occupancy are more accurate at coarse scales
 217 asymptotically: we can be surer that a species occupies Britain than it does some 1 km grid square
 218 therein. Our contribution has been to show that accuracy varies somewhat predictably along the
 219 spectrum from fine to coarse resolutions. Indeed, Meng's (2018) three-part decomposition of
 220 statistical error provides a clear theoretical framework within which analysts can consider quantities
 221 like the potential sampling bias and the sampling rate when deciding on the appropriate resolution at
 222 which to estimate occupancy. Coarsening the resolution will be particularly beneficial where sampling
 223 biases are likely to be large (e.g. when using citizen science data; Pescott et al., 2019, Stroh et al.,
 224 2023).

225 The Meng (2018) equation tells us that to increase the accuracy of estimates of species occupancy, we
 226 should work at the spatial resolution at which the sampling bias and the variability of occupancy in
 227 the landscape are smallest and at which the sampling rate is highest. Maximising the sampling rate is
 228 simplest in theory, because it always increases with resolution (practice of course introduces issues of
 229 resourcing and planning). The effect of resolution on the variability of occupancy in the landscape
 230 depends on the species' prevalence (i.e. \bar{Y}_N) at the finest resolution. If there is good reason to think
 231 that $\bar{Y}_N \geq 0.5$ —say, from an expert drawn range map—then coarsening the resolution will always
 232 reduce σ_Y . On the other hand, if there is good reason to think that \bar{Y}_N is truly low, then coarsening the
 233 resolution will increase σ_Y until the \bar{Y}_N reaches 0.5.

234 In our simulations, sampling bias was clearly lower at coarser resolutions (Fig. 2D), but this will not
235 be universally true. One minor thing to note is that we presented the average $\rho(R, Y)$ across 100
236 samples: for some of the individual samples, $\rho(R, Y)$ occasionally increased from one resolution to
237 the next. More general insight into how $\rho(R, Y)$ might scale with resolution in other situations can be
238 gleaned from the formula for Pearson's correlation coefficient. $\rho(R, Y)$ is the Pearson's correlation
239 between R and occupancy Y , which is to say their covariance divided by the product of their standard
240 deviations. We have already seen that coarsening the resolution of analysis increases the standard
241 deviation of Y σ_Y until $\bar{Y}_N \geq 0.5$, at which point further coarsening the resolution reduces it. The same
242 logic applies to the standard deviation of the R , which is also a binary variable. It follows that the
243 denominator in the formula for $\rho(R, Y)$, the product of the standard deviations of Y and R , increases
244 as the resolution is coarsened to the point where $\bar{Y}_N \geq 0.5$ and $P(R = 1) \geq 0.5$, at which point further
245 coarsening the resolution reduces it. For a given covariance between occupancy and R then,
246 coarsening the resolution of analysis will reduce $\rho(R, Y)$ where $\bar{Y}_N \geq 0.5$ and $P(R = 1) \geq 0.5$.
247 Further work is needed to understand how the covariance between occupancy and R will vary with
248 spatial resolution under different conditions.

249 As it is often time trends in species occupancy, rather than one-off estimates, that are of interest, it is
250 worth considering estimation error in this context. It is generally understood that time-varying
251 sampling bias (and therefore error) can confound true change in occupancy (Bowler et al., 2022), but
252 knowing how sampling bias changes over time is made difficult by the various sampling schemes and
253 analytical approaches that might be employed by researchers. The simplest scenario is where the
254 analyst estimates occupancy separately for multiple time-periods and calculates the differences
255 between them. If the sampling bias changes over time, then the estimated differences will be
256 erroneous. Another way to estimate time trends in occupancy is to restrict the analysis to the pool of
257 sites that were sampled at some point within the relevant timeframe and to predict (or impute) missing
258 values in each time-period (Boyd, August, et al., 2023; Isaac et al., 2014). Putting to one side the fact
259 that there are almost certain to be prediction errors, one ends up in a situation where the distribution of
260 R across sites is effectively time-invariant. Crucially, however, this does not mean that the sampling
261 bias will remain constant over time unless the distribution of Y across sites is also time-invariant (i.e.
262 the species' distribution does not change over time at the relevant scale). A similar scenario arises
263 when occupancy is estimated using unrepresentative monitoring data whose geographic distribution
264 does not change over time: for example, long-term monitoring of protected sites.

265 Understanding how the potential for confounding of error and true temporal change in occupancy
266 varies with spatial resolution is difficult, but the Meng equation provides several insights here too. For
267 example, working at coarser resolutions means less temporal variation in \bar{Y}_N (as colonisations and
268 local extinctions are less probable), which means less temporal variation in σ_Y . It is also likely to
269 mean less variation in $\rho(R, Y)$ —especially if occupancy is predicted across a fixed pool of sites in
270 each year, in which case the distribution of R is effectively constant over time (again, one must also
271 consider the fact that the predictions could be wrong at unsampled site/time-period combinations).
272 Reducing temporal variation in the quantities in eq. 1 will reduce temporal variation in error, which
273 should reduce the potential for confounding of error and true change in occupancy in many cases. An
274 obvious exception is where the per-period errors cancel each other out over long timeframes, in which
275 case they will not bias the estimated trend; however, it is not likely that biodiversity monitors will
276 know that they are in this situation—if the per period direction of error was known, then it could be
277 modelled. More elaborate simulations and theoretical work are needed to fully understand the effects
278 of spatial scale on error when estimating time trends in species occupancy.

279 The fact that error in estimates of species occupancy is likely to be lower at coarser spatial resolutions
280 sets up a trade-off between accuracy and “usefulness”. Estimates of species occupancy clearly have
281 the potential to be more useful at fine scales. For example, working at a finer resolution, at which
282 local extinctions and colonisations are more probable, means having a greater power to detect change.

283 (Of course, this argument supposes that the estimates are accurate or at least consistently inaccurate
284 over time. It also supposes that the power to detect change at some percentile is of primary interest,
285 which is not always true.) Other limitations of working at coarse resolutions are that occupancy is a
286 better surrogate for abundance, which is often of interest, and is often more relevant to policy at fine
287 scales (Kunin, 1998; Spake et al., 2022). When deciding on the appropriate resolution at which to
288 analyse their data, analysts must balance the need for accurate and useful estimates and remember that
289 an estimate will not be useful if it is completely wrong.

290 A good example of the potential for bias being balanced against the desire for finely resolved
291 estimates of species occupancy is found in the latest plant atlas of the Botanical Society of Britain and
292 Ireland (Stroh et al., 2023). The data were analysed at a 10×10 km scale—much coarser than the
293 1×1 km resolution used by others in the area (Boyd, August, et al., 2023)—and time-periods were
294 omitted, due to serious concerns about sampling biases affecting species data at finer scales across the
295 20th century. For example, both rarer and more challenging to identify taxa were more likely to be
296 reported at finer scales in the early part of the time series. Moreover, f was known to be far smaller at
297 smaller scales in these earlier periods (Pescott et al., 2019).

298 Like all simulations, ours are a simplification of reality, which might have implications for the wider
299 applicability of our results. We did not account for the fact that additional data tend to be available at
300 coarser resolutions; for example, digitised specimens may be resolved only to some vague locality,
301 and historic distribution data from species' Atlases tend to be more coarsely resolved than
302 contemporary data (Groom et al., 2018; Kunin et al., 2000; Pescott et al., 2019). These additional data
303 would increase the sampling rate f at coarse resolutions, which, as we have shown, would be likely to
304 increase the accuracy of sample-based estimates of mean occupancy. [Note that it is possible to
305 combine fine and coarse data using integrated distribution models and to draw inferences at the finer
306 scale (Pacifici et al., 2019). Whether the fact that data might be available solely at coarse scales for
307 historic time-periods, and at multiple scales for recent ones, will impact inference is an open
308 question.] Our assumption of perfect detection (i.e. no false absences) is also unrealistic, so it is worth
309 considering whether the prevalence of false absences is likely to be lower at fine or coarse resolutions.
310 On the one hand, if a coarse resolution is chosen when planning data collection, false absences might
311 be higher if the portions of the larger cells that are sampled are not suitable for the focal species
312 (Altwegg & Nichols, 2019). On the other, if the resolution is chosen at the analysis stage, coarsening
313 the spatial resolution increases the number of sampling events per site, so, all else being equal, it is
314 more likely that the focal species will be detected if it is present.

315 Rather than accepting false absences, it is common practice to try to correct them using some sort of
316 occupancy-detection model (MacKenzie et al., 2002; Royle, 2006). Coarsening the resolution of the
317 analysis risks violating the closure assumption of occupancy-detection models (Altwegg & Nichols,
318 2019; Jönsson et al., 2021), but it also increases the amount of repeat visits to the same site, which are
319 needed to estimate detectability and correct false absences. Interesting possibilities are that multi-scale
320 occupancy models (Mordecai et al., 2011), which relax the closure assumption, could be used and that
321 fine-scale sampling events could be used as spatial replicates to estimate detection probabilities and
322 correct false absences at coarser scales (Srivathsa et al., 2018). While failing to correct false absences
323 can make estimates of species occupancy worse, it is important to remember that successfully
324 correcting them only reduces error to its baseline level determined by sampling biases (Meng, 2018).

325 Coarsening the resolution of an analysis is one approach to counter some of the error introduced by
326 sampling biases, but there are alternatives. One is to estimate mean occupancy in the population using
327 a *weighted* sample mean, where the weights are equal to the inverse of the (possibly estimated)
328 sample inclusion probabilities (Boyd, Stewart, et al., 2023; Johnston et al., 2020). If successful,
329 weighting of this type brings the distribution of occupancy in the sample closer to its distribution in
330 the population and can be recast as a means to minimising $\rho(R, Y)$ (Meng, 2022). Several approaches

331 to estimating sampling weights for unstructured (i.e. nonprobability) samples, the principal type of
332 data used to estimate species occupancy, exist (Boyd, Stewart, et al., 2023; Elliott & Valliant, 2017).
333 Weighting is often more successful where available covariates explain larger portions of the variance
334 in sample inclusion (i.e. R) and the variable of interest (occupancy; (Collins et al., 2001), and it would
335 be useful to investigate how this scales with spatial resolution.

336 5 Conclusions

337 Analysts consider several factors when deciding on the appropriate resolution at which to estimate
338 species occupancy. Examples include the focal species' home range sizes (Wilson & Schmidt, 2015),
339 the scale at which they use the landscape more generally (Powney et al., 2019), the number of
340 replicate visits to the same site within closure periods (Outhwaite et al., 2019) and the resolution at
341 which the data were collected (Higa et al., 2015). We propose that analysts should also consider the
342 fact that estimates are likely to be more accurate at coarse resolutions, because a highly erroneous
343 finer-scale estimate is unlikely to be useful for most applications. The Meng (2018) equation provides
344 a theoretical framework in which accuracy and the desire for finely resolved information can be
345 balanced.

346 Acknowledgements

347 RJB and OLP were supported by the NERC Exploring the Frontiers award number NE/X010384/1
348 "Biodiversity indicators from nonprobability samples: Interdisciplinary learning for science and
349 society". All authors were supported by the NERC award number NE/R016429/1 as part of the UK-
350 SCAPE programme delivering National Capability.

351 Code availability

352 All code needed to fully reproduce our analysis is available at
353 <https://github.com/robboyd/biasVsResolution>.

354 References

- 355 Altwegg, R., & Nichols, J. D. (2019). Occupancy models for citizen-science data. *Methods in Ecology*
356 *and Evolution*, *10*(1), 8–21. <https://doi.org/10.1111/2041-210X.13090>
- 357 Azaele, S., Cornell, S. J., & Kunin, W. E. (2012). Downscaling species occupancy from coarse spatial
358 scales. *Ecological Applications*, *22*(3), 1004–1014. <https://doi.org/10.1890/11-0536.1>
- 359 Bowler, D. E., Callaghan, C. T., Bhandari, N., Henle, K., Barth, M. B., Koppitz, C., Klenke, R.,
360 Winter, M., Jansen, F., Bruelheide, H., & Bonn, A. (2022). Temporal trends in the spatial bias of
361 species occurrence records. *Ecography*. <https://doi.org/10.1111/ecog.06219>
- 362 Boyd, R. J., August, T., Cooke, R., Logie, M., Mancini, F., Powney, G., Roy, D., Turvey, K., & Isaac,
363 N. (2023). An operational workflow for producing periodic estimates of species occupancy at
364 large scales. *Biological Reviews*, *9*. <https://doi.org/10.32942/OSF.IO/2V7JP>
- 365 Boyd, R. J., Powney, G. D., & Pescott, O. L. (2023). We need to talk about nonprobability samples.
366 *Trends in Ecology & Evolution*, *38*(6), 521–531. <https://doi.org/10.1016/j.tree.2023.01.001>
- 367 Boyd, R. J., Stewart, G. B., & Pescott, O. L. (2023). Descriptive Inference using large,
368 unrepresentative nonprobability samples: An introduction for ecologists. *Ecology, forthcomin*.
369 <https://doi.org/10.1002/ecy.4214>
- 370 Collins, L. M., Schafer, J., & Kam, C. (2001). A Comparison of Restrictive Strategies in Modern
371 Missing Data Procedures. *Psychological Methods*, *6*(June). <https://doi.org/10.1037/1082-989X.6.4.330>

- 373 Dennis, E. B., Brereton, T. M., Morgan, B. J. T., Fox, R., Shortall, C. R., Prescott, T., & Foster, S.
374 (2019). Trends and indicators for quantifying moth abundance and occupancy in Scotland.
375 *Journal of Insect Conservation*, 23(2), 369–380. <https://doi.org/10.1007/s10841-019-00135-z>
- 376 Elliott, M. R., & Valliant, R. (2017). Inference for nonprobability samples. *Statistical Science*, 32(2),
377 249–264. <https://doi.org/10.1214/16-STS598>
- 378 Groom, Q. J., Marsh, C. J., Gavish, Y., & Kunin, W. E. (2018). How to predict fine resolution
379 occupancy from coarse occupancy data. *Methods in Ecology and Evolution*, 9(11), 2273–2284.
380 <https://doi.org/10.1111/2041-210X.13078>
- 381 Guélat, J., & Kéry, M. (2018). Effects of spatial autocorrelation and imperfect detection on species
382 distribution models. *Methods in Ecology and Evolution*, 9(6), 1614–1625.
383 <https://doi.org/10.1111/2041-210X.12983>
- 384 Hartley, S., & Kunin, W. E. (2003). Scale Dependency of Rarity, Extinction Risk, and Conservation
385 Priority. *Conservation Biology*, 17(6), 1559–1570. <https://doi.org/10.1111/j.1523-1739.2003.00015.x>
- 387 Higa, M., Yamaura, Y., Koizumi, I., Yabuhara, Y., Senzaki, M., & Ono, S. (2015). Mapping large-
388 scale bird distributions using occupancy models and citizen data with spatially biased sampling
389 effort. *Diversity and Distributions*, 21(1), 46–54. <https://doi.org/10.1111/ddi.12255>
- 390 Isaac, N. J. B., van Strien, A. J., August, T. A., de Zeeuw, M. P., & Roy, D. B. (2014). Statistics for
391 citizen science: Extracting signals of change from noisy ecological data. *Methods in Ecology
392 and Evolution*, 5(10), 1052–1060. <https://doi.org/10.1111/2041-210X.12254>
- 393 Johnston, A., Moran, N., Musgrove, A., Fink, D., & Baillie, S. R. (2020). Estimating species
394 distributions from spatially biased citizen science data. *Ecological Modelling*, 422(December
395 2019), 108927. <https://doi.org/10.1016/j.ecolmodel.2019.108927>
- 396 Jönsson, G. M., Broad, G. R., Sumner, S., & Isaac, N. J. B. (2021). A century of social wasp
397 occupancy trends from natural history collections: spatiotemporal resolutions have little effect
398 on model performance. *Insect Conservation and Diversity*, 14(5), 543–555.
399 <https://doi.org/10.1111/icad.12494>
- 400 Kéry, M., & Royle, J. A. (2016). *Applied hierarchical modelling in ecology: analysis of species
401 distribution, abundance and species richness in R and BUGS*. Academic press.
- 402 Kunin, W. E. (1998). Extrapolating species abundance across spatial scales. *Science*, 281(5382),
403 1513–1515. <https://doi.org/10.1126/science.281.5382.1513>
- 404 Kunin, W. E., Hartley, S., & Lennon, J. J. (2000). Scaling down: On the challenge of estimating
405 abundance from occurrence patterns. *American Naturalist*, 156(5), 560–566.
406 <https://doi.org/10.1086/303408>
- 407 MacKenzie, D. I., Nichols, J. D., Lachman, G. B., Droege, S., Royle, A. A., & Langtimm, C. A.
408 (2002). Estimating site occupancy rates when detection probabilities are less than one. *Ecology*,
409 83(8), 2248–2255. [https://doi.org/10.1890/0012-9658\(2002\)083\[2248:ESORWD\]2.0.CO;2](https://doi.org/10.1890/0012-9658(2002)083[2248:ESORWD]2.0.CO;2)
- 410 Meng, X.-L. (2018). Statistical paradises and paradoxes in big data (I): Law of large populations, big
411 data paradox, and the 2016 us presidential election. *Annals of Applied Statistics*, 12(2), 685–726.
412 <https://doi.org/10.1214/18-AOAS1161SF>

- 413 Meng, X.-L. (2022). Comments on the Wu (2022) paper by Xiao-Li Meng 1 : Miniaturizing data
414 defect correlation : A versatile strategy for handling non-probability samples. *Survey*
415 *Methodology*, 48(2), 1–22.
- 416 Mordecai, R. S., Mattsson, B. J., Tzilkowski, C. J., & Cooper, R. J. (2011). Addressing challenges
417 when studying mobile or episodic species: Hierarchical Bayes estimation of occupancy and use.
418 *Journal of Applied Ecology*, 48(1), 56–66. <https://doi.org/10.1111/j.1365-2664.2010.01921.x>
- 419 Outhwaite, C., Powney, G., August, T., Chandler, R., Rorke, S., Pescott, O. L., Harvey, M., Roy, H.
420 E., Fox, R., Roy, D. B., Alexander, K., Ball, S., Bantock, T., Barber, T., Beckmann, B. C., Cook,
421 T., Flanagan, J., Fowles, A., Hammond, P., ... Isaac, N. J. B. (2019). Annual estimates of
422 occupancy for bryophytes, lichens and invertebrates in the UK, 1970-2015. *Scientific Data*, 6(1),
423 259. <https://doi.org/10.1038/s41597-019-0269-1>
- 424 Pacifici, K., Reich, B. J., Miller, D. A. W., & Pease, B. S. (2019). Resolving misaligned spatial data
425 with integrated species distribution models. *Ecology*, 100(6), 1–15.
426 <https://doi.org/10.1002/ecy.2709>
- 427 Pescott, O. L., Humphrey, T. A., Stroh, P. A., & Walker, K. J. (2019). Temporal changes in
428 distributions and the species atlas: How can British and Irish plant data shoulder the inferential
429 burden? *British & Irish Botany*, 1(4), 250–282. <https://doi.org/10.33928/bib.2019.01.250>
- 430 Powney, G. D., Carvell, C., Edwards, M., Morris, R. K. A., Roy, H. E., Woodcock, B. A., & Isaac, N.
431 J. B. (2019). Widespread losses of pollinating insects in Britain. *Nature Communications*,
432 10(2019), 1–6. <https://doi.org/10.1038/s41467-019-08974-9>
- 433 Royle, J. A. (2006). Site occupancy models with heterogeneous detection probabilities. *Biometrics*,
434 62(1), 97–102. <https://doi.org/10.1111/j.1541-0420.2005.00439.x>
- 435 Spake, R., Barajas-Barbosa, M. P., Blowes, S. A., Bowler, D. E., Callaghan, C. T., Garbowski, M.,
436 Jurburg, S. D., Van Klink, R., Korell, L., Ladouceur, E., Rozzi, R., Viana, D. S., Xu, W. B., &
437 Chase, J. M. (2022). Detecting Thresholds of Ecological Change in the Anthropocene. *Annual*
438 *Review of Environment and Resources*, 47, 797–821. <https://doi.org/10.1146/annurev-environ-112420-015910>
- 440 Srivathsa, A., Puri, M., Kumar, N. S., Jathanna, D., & Karanth, K. U. (2018). Substituting space for
441 time: Empirical evaluation of spatial replication as a surrogate for temporal replication in
442 occupancy modelling. *Journal of Applied Ecology*, 55(2), 754–765.
443 <https://doi.org/10.1111/1365-2664.13005>
- 444 Stroh, P. A., Walker, K., Humphrey, T. A., Pescott, O. L., & Burkmar, R. (2023). *Plant Atlas 2020:*
445 *Mapping Changes in the Distribution of the British and Irish Flora*. Princeton Univ. Press.
- 446 Van Strien, A. J., Van Swaay, C. A. M., & Termaat, T. (2013). Opportunistic citizen science data of
447 animal species produce reliable estimates of distribution trends if analysed with occupancy
448 models. *Journal of Applied Ecology*, 50(6), 1450–1458. <https://doi.org/10.1111/1365-2664.12158>
- 450 Wilson, R. J., Thomas, C. D., Fox, R., Roy, D. B., & Kunin, W. E. (2004). Spatial patterns in species
451 distributions reveal biodiversity change. *Nature*, 432(7015), 393–396.
452 <https://doi.org/10.1038/nature03031>
- 453 Wilson, T., & Schmidt, J. H. (2015). Scale dependence in occupancy models: Implications for
454 estimating bear den distribution and abundance. *Ecosphere*, 6(9). <https://doi.org/10.1890/ES15-00250.1>
455

456

457