

Citizen scientists as butterfly predators: using foraging theory to understand individual recorder behaviour

Mingrui Li ^{a,*}, Robin J. Boyd ^b, Chloë Smith ^c, Richard Fox ^c, David Roy ^{a,b}, Jonathan Bennie ^d, Richard H. ffrench-Constant ^a

^a University of Exeter, Centre for Ecology and Conservation, Penryn Campus, Penryn, Cornwall, TR10 9FE, United Kingdom

^b UK Centre for Ecology and Hydrology, Benson Lane, Crowmarsh Gifford, Oxfordshire, OX10 8BB, UK

^c Butterfly Conservation, Manor Yard, East Lulworth, Wareham, Dorset, BH20 5QP, UK

^d University of Exeter, Centre for Geography and Environmental Science, Penryn Campus, Penryn, Cornwall, TR10 9FE, United Kingdom

ARTICLE INFO

Keywords:

Citizen science
Butterfly monitoring
iRecord
Sampling bias
Optimal foraging

ABSTRACT

Citizen science is increasingly important in the collection of biological data. However, to understand the broader utility of the growing number of citizen-derived records, we need to understand exactly how recorder behaviour affects the geographic distribution of records made. Here, we apply an optimal foraging model to citizen science data from the UK to determine how likely a recorder (predator) is to visit any given kilometre square and record a butterfly (prey). By defining the square with the highest density of an individual's records as their 'origin', we show that the probability of visiting a given site depends on its distance from the origin and the rarity-weighted species richness of the species thought to be present. This pattern of behaviour differs between recorders visiting more than or fewer than five squares, termed broad and narrow-range foragers. The model shows that recorder behaviour is driven, in part, by a simple trade-off between distance travelled and the rarity-weighted species richness. This collective behaviour helps explain over-recording by broad-ranging foragers in protected areas at distance and under-recording, by narrow-range foragers, in the wider countryside. It also implies that estimating parameters describing rare species' distributions (e.g. mean occupancy) will be challenging, since sample inclusion depends on occupancy itself. Mapping rare species' distributions should be simpler, since the sites at which they can be found tend to be well-sampled, but the same is unlikely to be true of common species, which also occupy areas that are unlikely to be sampled. More work is needed to understand how widely our results can be generalised beyond the UK and the dataset considered.

1. Introduction

Knowledge of species' distributions is fundamental to both ecological research and practical conservation efforts (Powney and Isaac, 2015). Activities such as biodiversity assessments, spatial conservation prioritisation, and environmental impact assessments all require information on how species are distributed geographically. While a map of a species' distribution may be purely empirical or obtained using a model, most ultimately derive from data on where particular species were observed (hereafter 'biological records'). It is therefore critical to recognise that biological records reflect the combined distribution of the species and the recorders themselves and that our knowledge of species distributions will therefore be biased by the spatial pattern of recording (Cretois et al., 2021; Geldmann et al., 2016; Hughes et al., 2021; Meyer

et al., 2016; Sicacha-Parada et al., 2021; Tiago et al., 2017).

Disentangling the distribution of the species and the recorders requires information on what motivated the recorders to collect data where they did (Boyd et al., 2025; Simmonds et al., 2020). While each recorder will have their own motivations, there are likely to be common factors that apply universally or at least approximately so. Since travelling incurs a cost (in terms of time and money), one obvious factor is the distance of the sampled location from where the recorder either lives or does most recording within the sampled area, which we term here as the recorder origin (Dennis and Thomas, 2000). Another is the perceived attractiveness of the location in terms of the species that might be seen there. Locally rare species, for example, might be perceived as high value targets (Dennis and Thomas, 2000; Bowler et al., 2022).

In many ways, the act of a recorder searching for 'interesting' species

* Corresponding author.

E-mail address: ml902@exeter.ac.uk (M. Li).

<https://doi.org/10.1016/j.ecolmodel.2025.111344>

Received 14 May 2025; Received in revised form 27 August 2025; Accepted 9 September 2025

Available online 16 September 2025

0304-3800/© 2025 The Authors. Published by Elsevier B.V. This is an open access article under the CC BY license (<http://creativecommons.org/licenses/by/4.0/>).

resembles a predator foraging for desirable prey. In both cases, individuals respond to attractants—such as rare species for recorders or high prey density for predators—and are deterred by other factors, like travel distance or competition. These conceptual similarities suggest that predator–prey theory may offer useful insights into patterns of biological recording.

In population and community models, site profitability is sometimes defined using a *functional response* (Boyd et al., 2020; Politikos et al., 2015). Basic forms describe how prey density affects ingestion rate (Holling, 1959), while more complex variants incorporate inhibitory effects, such as competitor density (e.g. Beddington et al. (1975), DeAngelis et al. (1975)). Since these models balance attractants and deterrents, they are a natural framework for understanding trade-offs in recorder behaviour.

Here, we use a simplified Beddington–DeAngelis-style functional response (Boyd et al., 2020) to examine why some sites are visited by citizen scientists while others are not. The model captures a trade-off between site desirability—defined as rarity-weighted species richness—and distance from the recorder’s origin. We focus on a subset of records submitted via the iRecord Butterflies app by citizen scientists, defined here as volunteer recorders making unstructured visits to sites of their choosing. Which species are recorded at sampled sites is not something that we consider.

2. Methods

2.1. Data sources

To assess how distance from the square at which most records were made and site desirability influence the probability that a given site was visited by a citizen scientist, we used data from two sources. The first dataset was derived from the iRecord Butterflies app, managed by the UK Centre for Ecology & Hydrology in partnership with Butterfly Conservation. This dataset consists of 51,045 records made by 1770 individual observers within Devon and Cornwall (Watsonian Vice Counties 1, 2, 3 and 4). Records spanned the period from 14 July 2011 to 1 April 2024. All records were included in our analysis, regardless of verification status, though most have been verified by local experts.

To estimate local species frequencies for the purpose of calculating site desirability (see below), we used a subset of the Butterflies for the New Millennium (BNM) national recording scheme, again focused on Devon and Cornwall. This subset comprises 785,582 records collected between 1796 and 31 December 2020. Only records that had been verified as correct or considered correct were retained for analysis. The BNM data serve as a more complete representation of local species occupancy, enabling us to derive measures of rarity-weighted species richness across 1 km × 1 km squares.

2.2. Square-specific metrics

We calculated three square-specific metrics to quantify site-level butterfly recording patterns: the number of individual records (N), the number of unique species recorded (S) and rarity-weighted species richness (R). The first two metrics, N and S, were derived from the iRecord Butterflies dataset. For R, we defined rarity-weighted species richness as the reciprocal of species frequency, where frequency was calculated as the number of 1 km squares in Devon and Cornwall where a given species was recorded in the BNM dataset. For a given square, the rarity-weighted species richness was the sum of these reciprocal values across all species recorded there from BNM dataset. This approach prioritises squares containing species that are locally rare, thereby providing a proxy for the desirability of a square to potential recorders.

2.3. Observer-specific metrics

To quantify recording behaviour at the level of individual observers,

we defined each observer’s origin as the 1 km square in which they recorded most total number of records. Using this as a spatial reference point, we calculated the cumulative number of records (N), species (S), and rarity-weighted species richness (R) as functions of distance from the origin square. For each observer, we determined the distance at which 50 % and 95 % of their total N, S, and R were accumulated. These values are referred to as 0.5 N and 0.95 N, 0.5 S and 0.95 S, and 0.5 R and 0.95 R, respectively. To illustrate this process, we constructed cumulative sampling curves for three known recorders based in Devon or Cornwall (Richard Fox, Marcus Rhodes, and Richard French-Constant) using UK-wide data from iRecord. These examples are intended for demonstration only and were not included in the core analysis, which focused exclusively on Devon and Cornwall.

2.3. Trade-off between distance and desirability

We used a simplified version of the Beddington–DeAngelis functional response (Beddington et al., 1975; DeAngelis et al., 1975) to model the trade-off between distance travelled from the observer origin square to each sampled square and the sum desirability of that square based on the rarity-weighted species richness. Importantly, we assume that recorder behaviour is primarily driven by this trade-off and that other potential influences such as accessibility and previous experience are negligible; the deterrent in our case is simply the distance from the recorder’s origin.

Let Y_{ij} be a binary variable indicating whether observer j visited site i , and let $P(Y_{ij})$ be the probability thereof. Then, we have

$$P(Y_{ij}) = \frac{R_i}{R_i + h_j + c_j D_{ij}},$$

where R_i is the rarity-weighted species richness of site i , D_{ij} is the distance between site i and observer j ’s origin, h_j represents observer j ’s baseline recording propensity (lower values imply greater propensities to record conditional on the effect distance) and c_j reflects the extent to which observer j is deterred by distance from their origin (higher values imply a reluctance to travel). The variables R_i and D_{ij} were derived as described above. Assuming the Y_{ij} s are independent Bernoulli trials, we used maximum likelihood estimation to obtain values for h_j and c_j .

We divided the individual recorders into two cohorts, narrow-range and broad-range foragers for simplicity, defined as those visiting fewer than or more than five squares beyond their origin square, respectively. We then fitted the predator prey model to each of the 1770 recorders individually and recorded their values of h and c . Finally, to examine the strength of the model in discriminating between sampled and non-sampled locations for the two different groups of recorders, we examined the area under the receiver-operator curve (AUC) for broad- and narrow-ranged foragers alike. The AUC takes a value between 0 and 1, with values closer to 1 indicating a better fit. To evaluate model performance for individual recorders, we also calculated both AUC and the correlation between predicted and observed values separately for each recorder, regardless of whether they were classified as broad- or narrow-ranged foragers.

We also reran our analysis using two alternative metrics of site attractiveness. The first was raw, unweighted species richness. The second was a subjective measure of cumulative species ‘charisma’, defined by one of the authors with extensive experience in butterfly recording in the region, and calculated as the sum of the charisma scores for all species present at the site.

3. Results

3.1. The number of records and distance travelled

Maps of either the total number of iRecord Butterflies records, the total number of species or rarity-weighted species richness present per 1

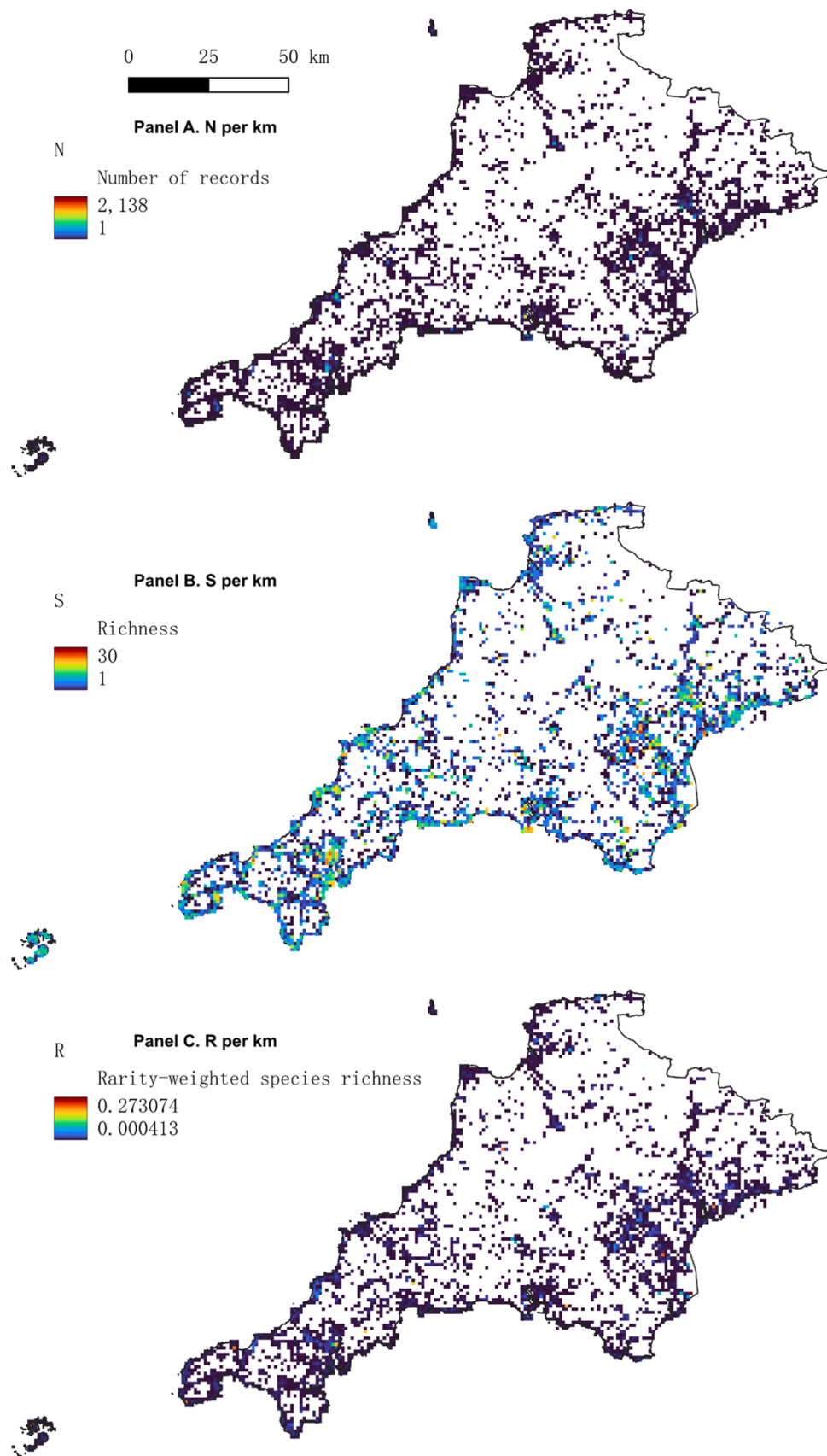


Fig. 1. Maps showing number of records (N), number of species (S) and the rarity-weighted species richness (R) for all species recorded from that kilometre square.

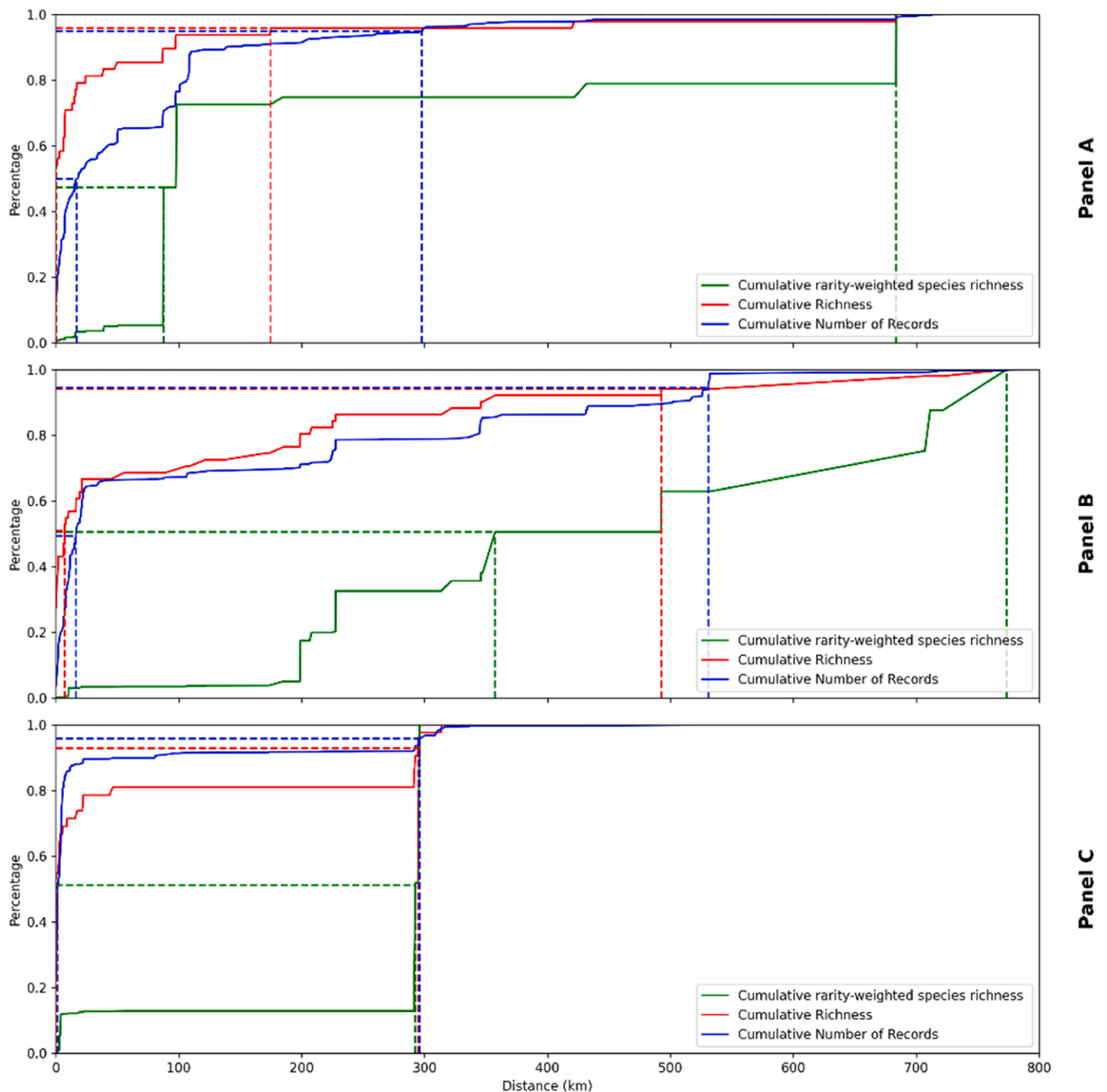


Fig. 2. Illustration of the cumulative sampling approach and metrics derived. Cumulative curves are shown for the number of records (N), species (S) and rarity-weighted species richness (R) recorded at distance (kilometres) from the observer origin square. Data are given for three known individual observers as three known recorders based in Devon or Cornwall. A. Richard Fox, B. Marcus Rhodes and C. Richard ffrench-Constant. All three observers are based in either Devon (RF) or Cornwall (MR and Rff-C) but records are drawn from across the whole of the UK. The dashed lines represent the distance travelled to record 50 or 95 % of each metric (0.5 or 0.95 N, 0.5 or 0.95 S and 0.5R or 0.95R). This same approach was then used to show the distance travelled by all observers within Devon and Cornwall only, see Fig. 3 and text for discussion.

km square (Fig. 1) show a clustered distribution. In fact, the strength of each signal around Penryn, which is the location of the Cornwall campus of the University of Exeter, suggests that the university may be a major local source of recorders and that these maps therefore represent the density and behaviour of individual recorders. As an example of how to quantify the distance travelled by individual recorders to make records, Fig. 2 shows the distance across the UK travelled to record 50 % or 95 % of records for three known recorders based in Devon or Cornwall. This simple illustration clearly shows that some recorders travel further to make their records than others, demonstrating that these metrics can be

used to quantify recorder behaviour in terms of distance travelled from their origin square. When these same 50–95 % metrics are used to examine the distribution of distance travelled by the combined population of all the 1770 observers in the iRecord Butterflies data for Devon and Cornwall (Fig. 3), we note that the statistics for most recorders are zero inflated in terms of distance. This shows that most recorders in Devon and Cornwall do not travel far from the kilometre square in which they do most of their recording.

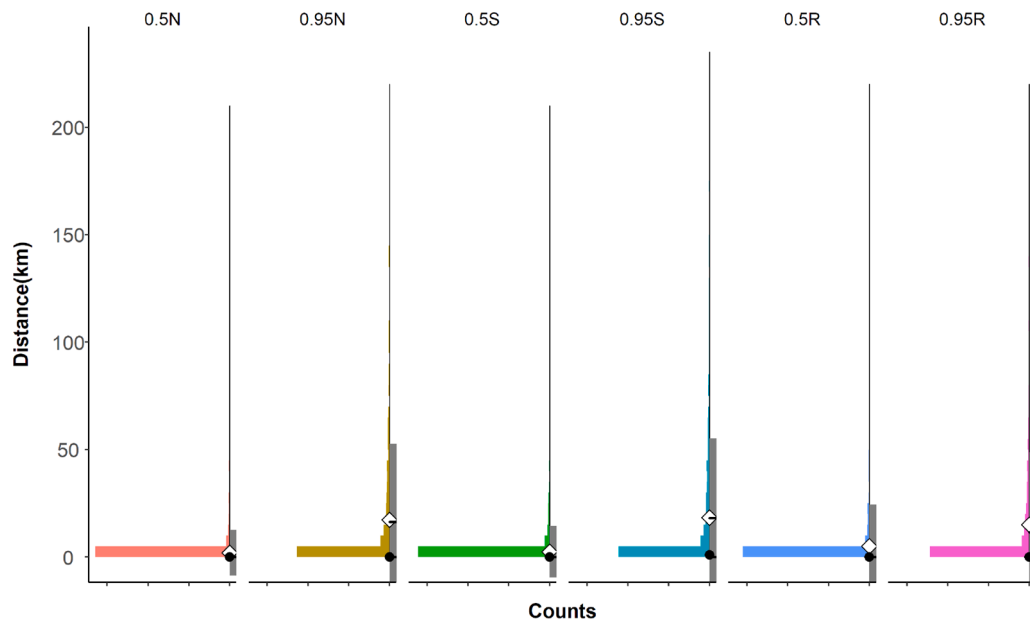


Fig. 3. Stack plots showing the distribution of 50 and 95 % values for the number of records (N, as 0.5 N and 0.95 N), number of species (S, as 0.5S and 0.95S) and rarity-weighted species richness (R, as 0.5R and 0.95R) for all 1700 recorders in the iRecord Butterflies dataset. The white diamond represents the mean and the vertical grey bar represents one standard deviation. The y-axis gives distance in km whereas counts are simply represented as ticks on the x-axis. Note that most recorders are zero inflated in terms of distance, suggesting most don't travel far from their origin square.

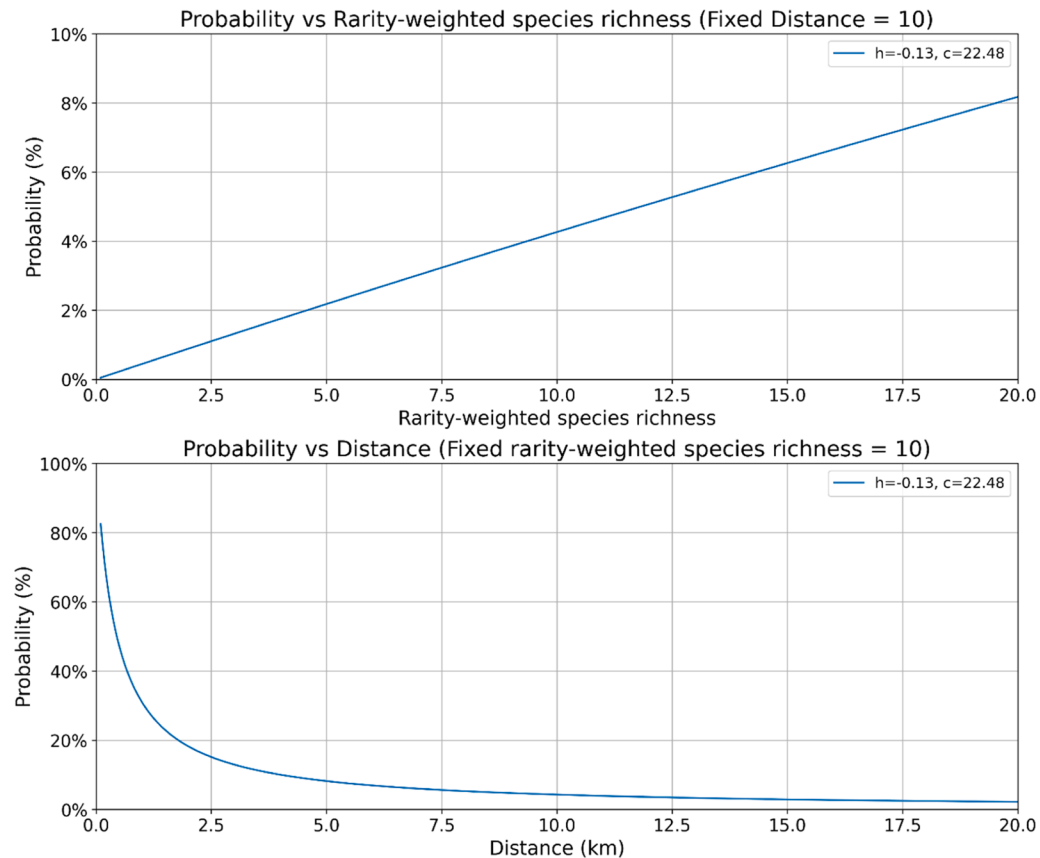


Fig. 4. Theoretical outputs from the predator-prey model showing the trade-off between rarity-weighted species richness (panel A) and distance from origin (panel B) in calculating the probability (%) that any given recorder will visit any given kilometre square in the sampled area. For panel A distance is fixed at 10 km, and the two constants h and c have values of -0.13 and 22.48 respectively. Whereas for panel B, rarity-weighted species richness is fixed at 10 and $h = -0.13$ and $c = 22.48$. Note that the two curves work in opposite directions and that the probability of any square being visited by any given observer is a therefore a trade-off of distance travelled and the desirability of the species target(s).

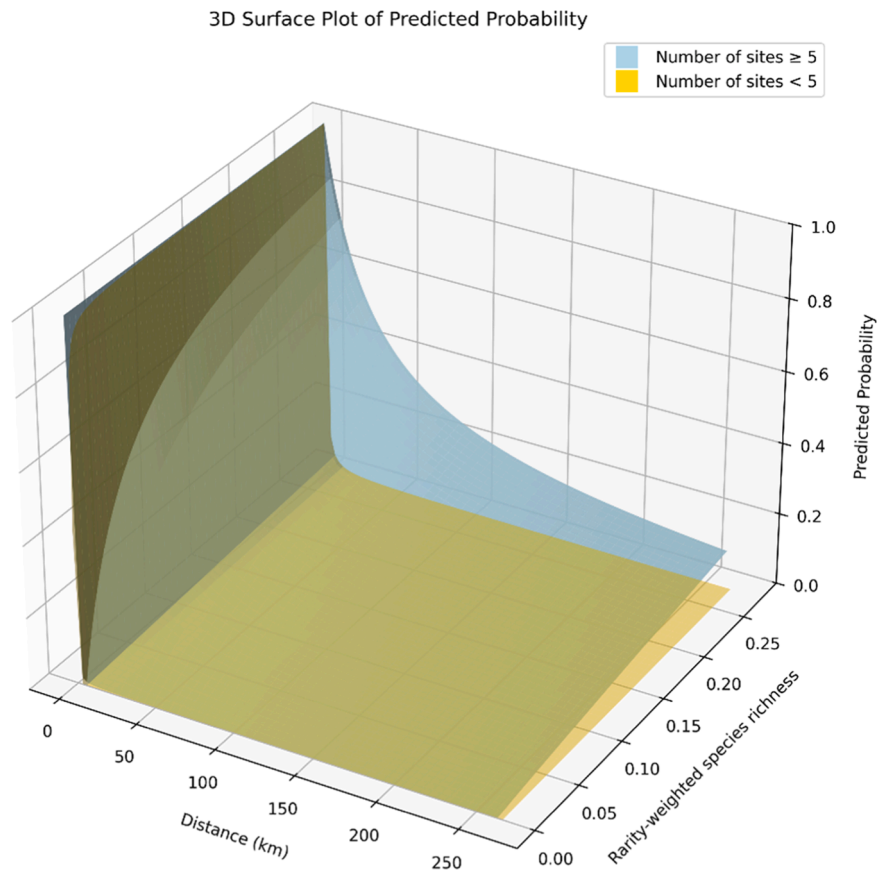


Fig. 5. Three-dimensional plot of the probability of any given square being visited by a recorder against distance from their likely origin and rarity-weighted species richness known from that square. The yellow curve represents recorders visiting <5 squares (termed narrow-range foragers) and the blue curve represents those visiting >5 squares beyond their origin (broad-range foragers). Note that experts are prepared to travel greater distances to visit any given square.

3.2. The trade-off between distance and desirability

To illustrate the likely trade-off between distance travelled from the square with the highest density of records and the desirability of species sought by the individual observer, we fitted the functional response to the iRecord Butterflies data. As an illustration of the model, by fixing the rarity-weighted species richness (perceived desirability) of species present at a relative value of 10, we show how increasing distance from the origin square decreases the probability that any given square will be sampled (Fig. 4A). Likewise, by fixing the distance travelled by an observer to an arbitrary 10 km we show how increasing the rarity-weighted species richness (desirability) increases the probability that any given square will be visited by the observer (Fig. 4B). The fit of the model in terms of AUC and $\text{corr}(Y, \hat{P}[Y])$ (see supplementary material) suggests that distance from their origin square and the rarity-weighted species richness explain a large proportion of the variation in where citizen scientists recorded.

3.3. Individual recorder behaviour

Plotting distance from observer origin against desirability of species present against the probability that any given square is visited (Fig. 5), we can see that narrow-range foragers record very close to, or within, their origin square. In contrast, broad-range foragers travel further (up to 250 km) to record species with higher perceived rarity. This effect can be quantified by examining the relative distribution of the constants h and c between the two different groups of recorders (Fig. 6). Broad-range foragers have lower values of c , showing that they are less deterred by travelling greater distances to record. In terms of h , narrow-range foragers sample fewer grid squares, and therefore appear less motivated in

terms of spatial coverage. However, after accounting for the deterrent effect of distance, the model may still assign them lower h values (Fig. 6). In this sense, h functions as a normalizing constant. The Receiver Operating Characteristic (ROC) analysis gives area under the curve (AUC) values of 0.97 for narrow-range foragers and 0.89 for broad-range foragers using the rarity-weighted species richness as the basis (Fig. 7).

We also tested the alternative metrics. When using raw, unweighted species richness, the model achieved an AUC of 0.88 for narrow-range foragers and 0.94 for broad-range foragers. Using a subjective measure of cumulative species charisma, defined by an experienced butterfly recorder as the sum of desirability scores for all species present, the AUC values were 0.84 and 0.93 respectively. In comparison, the model based on rarity-weighted species richness showed higher AUC values for both forager types, indicating better performance than either alternative metric (see supplementary figures).

4. Discussion

We applied a functional response, previously used to describe the dependence of predator consumption rate on competitor and prey density, to understand the geographic distribution of butterfly records made by citizen scientists in the UK. The model is built on first principles and supposes that recorders trade-off the perceived attractiveness of a site (in terms of rare species present) with its distance from their origin when deciding whether it is worthwhile visiting. We fitted the models to data for each recorder and on aggregate for two groups of recorders: 'broad-ranged' foragers, who have visited more than five 1 km squares, and 'narrow-ranged' foragers, who have visited fewer than five 1 km squares. The models explained a surprisingly large amount of the

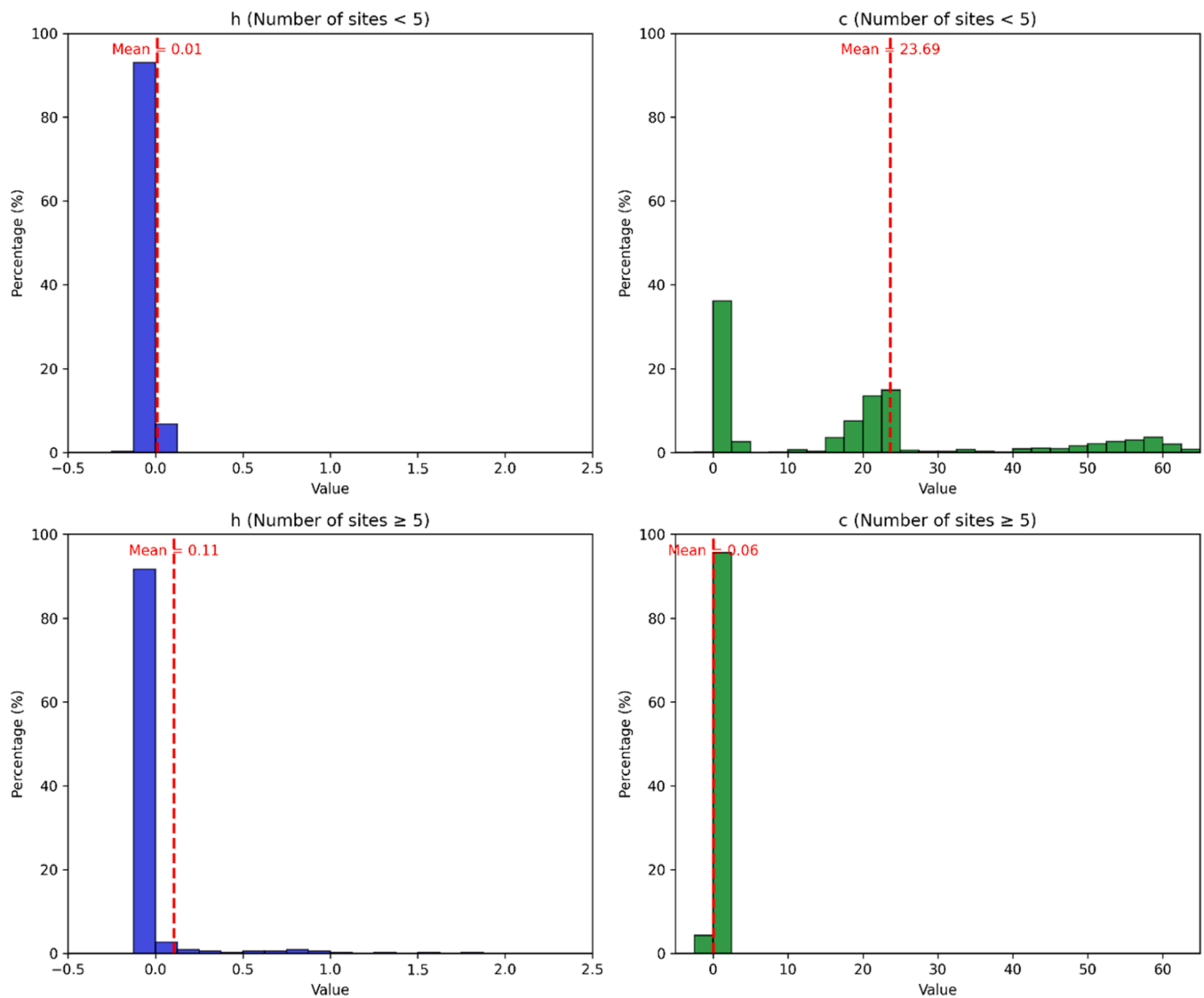


Fig. 6. The distribution of the two constants h and c from the predator-prey model for all recorders in the dataset with recorders again divided into broad-range foragers and narrow-range foragers. The constant c represents the strength of the deterrence of distance on recorders and h is the baseline probability that any given square will be sampled. Note that broad-range foragers have lower values of c , showing that they are less deterred by distance and are therefore more likely to sample squares at distance. The constant h is addressed in the main text.

variation in which sites were sampled, although the models for some individual recorders did not perform well (supplementary material).

It is not surprising that some of the models for individual recorders did not perform well. Some recorders will be motivated by factors not included in our models, such as the desire to record species that they have not previously observed (Goldstein and Stoudt, 2025). And for those recorders not deterred by the need to travel and the absence of rare species, the models are naturally less able to discriminate between sampled and non-sampled sites. Indeed, August et al. (2020) showed that many users of the iRecord Butterflies app do not focus on rare species and collect data over vast geographic areas. The models might also perform poorly for transient recorders (e.g. holidaymakers), whose origins might change over time; transient recorders often dominate citizen science datasets (Dimson and Gillespie, 2023), and Devon and Cornwall are popular tourist destinations. Nevertheless, it is clear that, on aggregate, distance from recorders' origins and the rarity-weighted species richness present are two key determinants of which sites are featured in the iRecord Butterflies dataset.

We use rarity-weighted species richness as the measure of site attractiveness in our model. The estimation of parameters h and c was

conducted under the assumption that recorders have imperfect knowledge of this index, as such imperfect knowledge is implicitly embedded in the spatial structure of the data and influences the locations that recorders choose to visit. As a result, the parameter estimates already reflect this cognitive bias to some extent, and this bias, together with the distribution of rarity-weighted species richness, shapes the way observers perceive site attractiveness. The strong fit of the model to the observed data further supports the validity of this metric.

The model parameters (h and c) have different values for different types of recorders. This may lie in variation in recorder skill and field-craft (Kühn et al., 2024). Individuals with greater taxonomic knowledge and field experience are likely to detect and identify a wider range of species while casual recorders or those with limited identification ability may tend to focus on a narrower subset of more recognisable or easily observed species (Kühn et al., 2024). Because more skilled recorders are able to identify a wider range of species, including rare ones, they may be more likely to visit sites with higher species diversity or those known to support particular target species (Bowler et al., 2022).

Our results have two major implications, the first being that estimating parameters describing rare species' distributions is

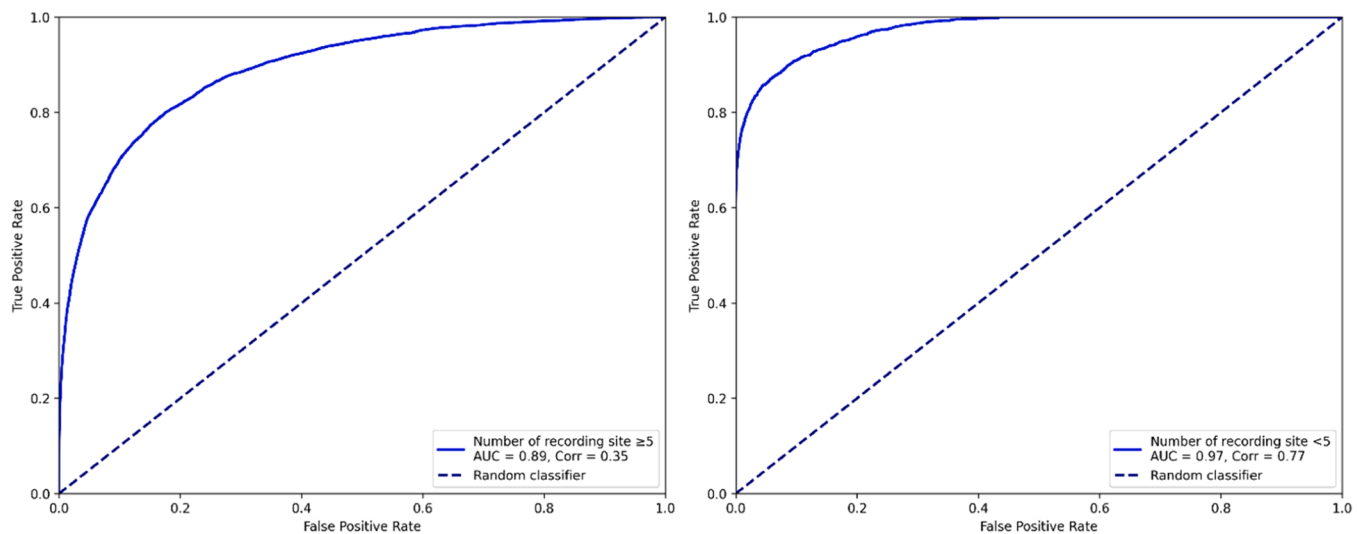


Fig. 7. Receiver Operating Characteristic (ROC) curves evaluating the ability of a cumulative rarity-based model to predict butterfly presence, separated by observer experience. The left panel shows results for broad-range foragers (observers with ≥ 5 recording sites), and the right panel for narrow-range foragers (< 5 recording sites). Differences in Area Under the Curve (AUC) values reflect variation in model performance between the two observer groups (see text for discussion).

fundamentally compromised by non-ignorable sampling. According to our theory, sites occupied by rare species are specifically targeted (Fig. 5), which means the outcome variable (occupancy of a rare species) at a given site directly affects the probability that is sampled. This is a classic case of Missing Not At Random (MNAR): the process that determines whether a site appears in the dataset depends on the outcome (Bowler et al., 2025; Little and Rubin, 2019). Under MNAR, statistical corrections are not possible without very strong assumptions, which are usually untestable. The silver lining is that, because sites at which rare species can be found are often well-sampled, mapping their distributions may still be feasible—even if unbiased estimation of descriptive parameters (e.g. mean occupancy) is not.

The second major implication of our results is that thoroughly mapping commoner species is likely to be difficult. The probability that any given site was sampled was determined by its proximity to recorders' origin squares and rarity-weighted species richness. Common species are likely to occupy sites not occupied by rare species and far from where people do most of their recording. Hence, many areas that appear to be unoccupied by common species will in fact just not be sampled.

Of course, rather than mapping the distribution of records or estimating parameters describing species' distributions naively, one might apply some sort of statistical correction. Most correction procedures involve controlling for variables that induced the bias in the first place. Examples include species distribution modelling (assuming the biasing variables are included as covariates), quasi-randomisation (i.e. propensity score weighting) and poststratification (Boyd et al., 2024).

Our findings have important implications for the effectiveness of statistical correction procedures that control for biasing variables. We showed that rare species' populations are preferentially sampled, which is a classic case of non-sampled sites being Missing Not At Random. For these species, controlling for biasing variables will not be sufficient to eliminate bias, and alternative approaches that require stronger assumptions will be needed (see e.g. Bailey (2023)). For commoner species, our findings imply that one should control for rare species' distributions where these can be approximated (e.g. from an Atlas, a field guide or based on knowledge of habitat requirements).

Alternative approaches to mitigating spatial sampling biases that do not require controlling for biasing variables might also be considered. When additional data sources, with different or even no biases are available, these might be integrated with the original sample. Alternatively, where one has some control over where new data are collected,

adaptive sampling of underrepresented habitats might be considered (Callaghan et al., 2019).

Our model, of course, has some limitations. While distance from the origin was a useful proxy for recording effort, it did not fully account for spatial and temporal variability in effort. Factors such as accessibility by public transport (Mair and Ruete, 2016; Sicacha-Parada et al., 2021; Mandeville et al., 2022) might also shape where and how recording would take place. Secondly, our study used distance from the origin square as a proxy for competition among recorders. However, this approach does not fully capture the complexity of recorder behaviour (Isaac and Pocock, 2015) and multiple competing biases between recorders may affect model predictions (Bowler et al., 2022). Thirdly, the relative dominance of a particular forager type within a dataset might significantly influence model outcomes (August et al., 2020). If there are loads of narrow ranged foragers, then the dataset will be dominated by areas with high population density; if the dataset is dominated by broad ranged foragers, the dataset will be dominated by places where there are rare species. Lastly, our model was based on opportunistic monitoring, where observers were free to choose where, when and what to report (Soroye et al., 2018). The current model would not be applicable for volunteers following systematic protocols that direct location of recording, such as those used in the UK Butterfly Monitoring Scheme. Therefore, such limitations and bias might indicate competitive dynamics among recorders and should be explored in the future.

In future work, we also intend to extend the current framework to account for the selection of species recorded once a site is visited. This component of recorder behaviour is likely shaped by taxonomic preferences (Goldstein et al., 2024), preference for recording previously unseen species (Soroye et al., 2018), sampling completeness (Sánchez-Fernández et al., 2021) and an emphasis on rare species (Habel et al., 2025). Capturing these dimensions would require the development of a new model that goes beyond spatial decision-making and describes species level choices made by observers. This represents an important next step toward a more comprehensive understanding of bias in opportunistic data collection.

A further promising direction for future work is to test the model's transferability to a different, yet comparable, taxonomic group. In Britain, the Cerambycidae is a family roughly similar in size to butterflies and with many species readily identifiable from photographs (Alexander, 2019). Importantly, occurrence records for Cerambycidae are openly available via iRecord. Utilizing these opportunistic datasets allows us to evaluate observer behavior patterns and model performance

beyond butterflies, thus strengthening the generality and applicability of our framework.

It is clear that the geographic distribution of records in the iRecord Butterflies dataset is influenced by the trade-off between the perceived desirability of a given 1 km square and its distance from recorders' origins. It follows that non-sampled squares are 'Missing Not At Random' and that drawing accurate inferences about species' distributions from the dataset will be challenging. Further work is needed to understand whether this fundamental trade-off applies to other citizen science datasets around the world.

File 4 – Supplementary material (hosted on Figshare)

The term "Supplementary material" in our repository refers to an internal folder within our Figshare archive, not separate journal-hosted supplementary files. The dataset and associated scripts are archived at Figshare DOI: 10.6084/m9.figshare.28847306.

CRediT authorship contribution statement

Mingrui Li: Writing – review & editing, Methodology, Formal analysis, Data curation. **Robin J. Boyd:** Writing – review & editing, Methodology, Conceptualization. **Chloë Smith:** Writing – review & editing, Resources. **Richard Fox:** Writing – review & editing, Methodology. **David Roy:** Writing – review & editing. **Jonathan Bennie:** Writing – review & editing, Methodology. **Richard H. ffrench-Constant:** Writing – review & editing, Writing – original draft, Supervision, Conceptualization.

Declaration of competing interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

Acknowledgements

We thank all present and past county recorders for Devon and Cornwall for their verification of the butterfly records and to all of the recorders who contributed sightings. We also thank Butterfly Conservation and the UK Centre for Ecology & Hydrology for provision of the iRecord Butterflies data set and Butterfly Conservation for the extract of Butterflies for the New Millennium recording scheme data. We are grateful to two anonymous reviewers for their comments on a previous version of this manuscript. RF was supported by the Heather Corrie Fund from Butterfly Conservation. For the purpose of open access, the author has applied a Creative Commons Attribution (CC BY) licence to any Author Accepted Manuscript version arising from this submission.

Data availability

All research data have been shared on FigShare at <https://doi.org/10.6084/m9.figshare.28847306>

References

- Alexander, K.N.A., 2019. A Review of the Status of the Beetles of Great Britain: Longhorn Beetles (Cerambycidae). Natural England, London.
- August, T., Fox, R., Roy, D.B., Pocock, M.J.O., 2020. Data-derived metrics describing the behaviour of field-based citizen scientists provide insights for project design and modelling bias. *Sci. Rep.* 10, 11009.
- Bailey, M.A., 2023. A new paradigm for polling. *Harv. Data Sci. Rev.* 5, 1–27.
- Beddington, J., Free, C., Lawton, J., 1975. Dynamic complexity in predator-prey models framed in difference equations. *Nature* 255, 58–60.
- Bowler, D.E., Bhandari, N., Repke, L., Beuthner, C., Callaghan, C.T., Eichenberg, D., Henle, K., Klenke, R., Richter, A., Jansen, F., Bruelheide, H., Bonn, A., 2022. Decision-making of citizen scientists when recording species observations. *Sci. Rep.* 12, 11069.
- Bowler, D.E., Boyd, R.J., Callaghan, C.T., Robinson, R.A., Isaac, N.J.B., Pocock, M.J.O., 2025. Treating gaps and biases in biodiversity data as a missing data problem. *Biolog. Rev.* 100, 50–67.
- Boyd, R.J., Botham, M., Dennis, E., Fox, R., Harrower, C., Middlebrook, I., Roy, D.B., Pescott, O.L., 2025. Using causal diagrams and superpopulation models to correct geographic biases in biodiversity monitoring data. *Methods Ecol. Evol.* 16, 332–344.
- Boyd, R.J., Sibly, R., Hyder, K., Walker, N., Thorpe, R., Roy, S., 2020. Simulating the summer feeding distribution of Northeast Atlantic mackerel with a mechanistic individual-based model. *Prog. Oceanogr.* 183, 102299.
- Boyd, R.J., Stewart, G.B., Pescott, O.L., 2024. Descriptive inference using large, unrepresentative nonprobability samples: an introduction for ecologists. *Ecology* 105, e4214.
- Callaghan, C.T., Rowley, J.J., Cornwell, W.K., Poore, A.G., Major, R.E., 2019. Improving big citizen science data: moving beyond haphazard sampling. *PLoS Biol.* 17, e3000357.
- Cretois, B., Simmonds, E.G., Linnell, J.D., Van Moorter, B., Rolandsen, C.M., Solberg, E. J., Strand, O., Gundersen, V., Roer, O., Rid, J.K., 2021. Identifying and correcting spatial bias in opportunistic citizen science data for wild ungulates in Norway. *Ecol. Evol.* 11, 15191–15204.
- Deangelis, D.L., Goldstein, R., O'Neill, R.V., 1975. A model for tropic interaction. *Ecology* 56, 881–892.
- Dennis, M., Roger, Thomas, C., 2000. Bias in butterfly distribution maps: the influence of hot spots and recorder's home range. *J. Insect Conserv.* 4, 73–77.
- Dimson, M., Gillespie, T.W., 2023. Who, where, when: observer behavior influences spatial and temporal patterns of iNaturalist participation. *Appl. Geogr.* 153, 102916.
- Geldmann, J., Heilmann Clausen, J., Holm, T.E., Levinsky, I., Markussen, B., Olsen, K., Rahbek, C., Tittrop, A.P., 2016. Exploring four recording schemes with different proficiency requirements. *Diver. Distribut.* 22, 1139–1149.
- Goldstein, B., Stoudt, S., Lewthwaite, J., Shirey, V., Mendoza, E., Guzman, M., 2024. Logistical and preference bias in participatory science butterfly data. *Frontiers in Ecology and the Environment*, p. 22.
- Goldstein, B.R., Stoudt, S., 2025. Evidence of novelty and specialization behavior in participatory science reporting. *Oikos*, e10938 n/a.
- Habel, J.C., Schmitt, T., Huemer, P., Rüdiger, J., Gros, P., Ulrich, W., 2025. Selective observation causes differences in citizen science butterfly data. *Basic Appl. Ecol.* 87, 46–54.
- Holling, C.S., 1959. Some characteristics of simple types of predation and parasitism. *Can. Entomol.* 91, 385–398.
- Hughes, A.C., Orr, M.C., Ma, K., Costello, M.J., Waller, J., Provoost, P., Yang, Q., Zhu, C., Qiao, H., 2021. Sampling biases shape our view of the natural world. *Ecography* 44, 1259–1269.
- Isaac, N.J.B., Pocock, M.J.O., 2015. Bias and information in biological records. *Biolog. J. Linnean Soc.* 115, 522–531.
- Kühn, E., Harpke, A., Schmitt, T., Settele, J., Kühn, I., 2024. Counting butterflies—Are old-fashioned ways of recording data obsolete? *J. Insect Conserv.* 28, 577–588.
- Little, R.J., Rubin, D.B., 2019. Statistical Analysis With Missing Data. John Wiley & Sons.
- Mair, L., Ruete, A., 2016. Explaining spatial variation in the recording effort of citizen science data across multiple taxa. *PLoS One* 11, e0147796.
- Mandeville, C.P., Nilsen, E.B., Finstad, A.G., 2022. Spatial distribution of biodiversity citizen science in a natural area depends on area accessibility and differs from other recreational area use. *Ecolog. Solut. Evid.* 3, e12185.
- Meyer, C., Weigelt, P., Kreft, H., 2016. Multidimensional biases, gaps and uncertainties in global plant occurrence information. *Ecol. Lett.* 19, 992–1006.
- Politikos, D.V., Huret, M., Petitgas, P., 2015. A coupled movement and bioenergetics model to explore the spawning migration of anchovy in the Bay of Biscay. *Ecol. Modell.* 313, 212–222.
- Powney, G.D., Isaac, N.J.B., 2015. Beyond maps: a review of the applications of biological records. *Biolog. J. Linnean Soc.* 115, 532–542.
- Sánchez-Fernández, D., Fox, R., Dennis, R.L., Lobo, J.M., 2021. How complete are insect inventories? An assessment of the British butterfly database highlighting the influence of dynamic distribution shifts on sampling completeness. *Biodivers. Conserv.* 30, 889–902.
- Siccha-Parada, J., Steinsland, I., Cretois, B., Borgelt, J., 2021. Accounting for spatial varying sampling effort due to accessibility in Citizen Science data: a case study of moose in Norway. *Spat. Stat.* 42, 100446.
- Simmonds, E.G., Jarvis, S.G., Henrys, P.A., Isaac, N.J.B., O'hara, R.B., 2020. Is more data always better? A simulation study of benefits and limitations of integrated distribution models. *Ecography* 43, 1413–1422.
- Soroye, P., Ahmed, N., Kerr, J.T., 2018. Opportunistic citizen science data transform understanding of species distributions, phenology, and diversity gradients for global change research. *Glob. Chang. Biol.* 24, 5281–5291.
- Tiago, P.C., Ceia-Hasse, A., Marques, T.A., Capinha, C.S., Pereira, H.M., 2017. Spatial distribution of citizen science casuistic observations for different taxonomic groups. *Sci. Rep.* 7, 12832.