

## Realising the promise of large data and complex models

McCrea, Rachel; King, Ruth; Graham, Laura; Börger, Luca

DOI:

[10.1111/2041-210X.14050](https://doi.org/10.1111/2041-210X.14050)

License:

Creative Commons: Attribution (CC BY)

*Document Version*

Publisher's PDF, also known as Version of record

*Citation for published version (Harvard):*

McCrea, R, King, R, Graham, L & Börger, L 2023, 'Realising the promise of large data and complex models', *Methods in Ecology and Evolution*, vol. 14, no. 1, pp. 4-11. <https://doi.org/10.1111/2041-210X.14050>

[Link to publication on Research at Birmingham portal](#)

### General rights

Unless a licence is specified above, all rights (including copyright and moral rights) in this document are retained by the authors and/or the copyright holders. The express permission of the copyright holder must be obtained for any use of this material other than for purposes permitted by law.

- Users may freely distribute the URL that is used to identify this publication.
- Users may download and/or print one copy of the publication from the University of Birmingham research portal for the purpose of private study or non-commercial research.
- User may use extracts from the document in line with the concept of 'fair dealing' under the Copyright, Designs and Patents Act 1988 (?)
- Users may not further distribute the material nor use it for the purposes of commercial gain.

Where a licence is displayed above, please note the terms and conditions of the licence govern your use of this document.

When citing, please reference the published version.

### Take down policy

While the University of Birmingham exercises care and attention in making items available there are rare occasions when an item has been uploaded in error or has been deemed to be commercially or otherwise sensitive.

If you believe that this is the case for this document, please contact [UBIRA@lists.bham.ac.uk](mailto:UBIRA@lists.bham.ac.uk) providing details and we will remove access to the work immediately and investigate.

## EDITORIAL

## Realising the Promise of Large Data and Complex Models

## Realising the promise of large data and complex models

Rachel McCrea<sup>1</sup>  | Ruth King<sup>2</sup>  | Laura Graham<sup>3,4</sup>  | Luca Börger<sup>5,6</sup> 

<sup>1</sup>Department of Mathematics and Statistics, Lancaster University, Lancaster, UK; <sup>2</sup>School of Mathematics and Maxwell Institute for Mathematical Sciences, University of Edinburgh, Edinburgh, UK; <sup>3</sup>Geography, Earth & Environmental Sciences, University of Birmingham, Birmingham, UK; <sup>4</sup>Biodiversity, Ecology & Conservation Group, International Institute for Applied Systems Analysis, Vienna, Austria; <sup>5</sup>Department of Biosciences, Swansea University, Swansea, UK and <sup>6</sup>Centre for Biomathematics, Swansea University, Swansea, UK

## Correspondence

Rachel McCrea

Email: [r.mccrea@lancaster.ac.uk](mailto:r.mccrea@lancaster.ac.uk)

## Funding information

Natural Environment Research Council, Grant/Award Number: NE/T009373/1; EPSRC, Grant/Award Number: EP/S020470/1

Handling Editor: Aaron Ellison

## 1 | MOTIVATION FOR THIS SPECIAL FEATURE

In an era of rapid change, ecologists are increasingly asked to provide answers to big, urgent questions of global concern (Solé & Levin, 2022; Sutherland et al., 2013; Yates et al., 2018). Concurrently, technological advances allow ecological data to be collected at increasingly high resolutions (e.g. temporal and/or spatial scales), leading to both new types of data and larger datasets becoming available (Farley et al., 2018). These data provide the opportunity to investigate new, and even previously unanswerable, questions, including those concerning animal movements (Nathan et al., 2022) and those addressing conservation and sustainability issues (Runting et al., 2022). Increasingly, realistic models need to be developed and fitted to these data (Fer et al., 2018), pushing the boundaries of the type and intricacy of questions that can be explored (Niu et al., 2020). However, big data and big models can lead to big troubles across multiple aspects, from storing and processing the data to fitting of complex models to data and interpreting the output.

Close collaborations between ecologists, statisticians, mathematical modellers, computer scientists and other disciplines offer exciting ways forward to solve these problems, leading to mutually beneficial advancements. For example, computer scientists can aid in the efficient storage and extraction of data, and the development of new algorithms; statisticians can help and guide ecologists in the analysis of data, fitting complex models to the data via efficient

computational algorithms and propagating or quantifying uncertainties throughout the process; mathematicians can ensure models are constructed in the most suitable fashion for the specific questions asked and demonstrate suitable properties (such as realistic territorial ranges or population predictions); and ecologists can guide mathematical scientists on the biological characteristics of the systems studied and ecological interpretation of the corresponding results, thus informing future models and influencing policy decisions. The need to answer important ecological questions is unprecedented, with declines in biodiversity and ecosystem services which will impact our ability to meet Sustainable Development Goals (Reyers & Selig, 2020), and it is through interdisciplinary collaborations that the biggest steps forward can be made.

Data analysis challenges arise across the full data analytic pipeline, including processing and visualising the data, developing ecologically relevant and interpretable models to fit to the data, adapting the associated algorithms to fit models to data efficiently and obtaining meaningful interpretations of the output. In practice, there are often many trade-offs between these different aspects due to the challenges that arise during the data analysis pipeline. For example, within the initial processing of the data, decisions may need to be made regarding cleaning the data (e.g. to remove recorded data errors) or the summarised form of the processed data to report (e.g. the temporal and/or spatial scale). This itself can be challenging and there will often be uncertainty within the process, leading to potential new errors being introduced. The

This is an open access article under the terms of the [Creative Commons Attribution](https://creativecommons.org/licenses/by/4.0/) License, which permits use, distribution and reproduction in any medium, provided the original work is properly cited.

© 2023 The Authors. *Methods in Ecology and Evolution* published by John Wiley & Sons Ltd on behalf of British Ecological Society.

decisions made will typically impact the model fitted to these data. For example, for motion-sensor camera trap data, there may be a trade-off between the level of initial data processing (i.e. the level of advanced tools used for uniquely identifying individuals via, e.g. machine learning techniques) and associated models that may be fitted to incorporate the amount of uncertainty in the preprocessed data (e.g. from assuming no error in the matches; to incorporating matching uncertainty; to allowing for both marked and unmarked individuals). Alternatively, complex models often require computationally intensive algorithms for them to be fitted to the data, which may not scale as datasets increase in size. This may lead to the consideration of a simpler model that can be more easily fitted, thus reducing the level of fine-detail that may be extracted from the data; or adaptations to the model-fitting process such as using some form of approximate model-fitting approach that aims to be robust to the approximations used, but potentially could lead to biased parameter estimates.

This Special Feature provides a combination of review papers and scientific articles that address one or more of the challenges of modern day analyses of large and/or complex ecological data. Echoing the challenges facing the discipline, we present these in the natural statistical cycle, starting with the challenges of new types of data, to the limitations of statistical models and associated algorithms (and computer packages) used to fit the models to the data to the interpretation and presentation of the corresponding model outputs.

## 2 | BROAD THEMES

We consider each of the themes identified in turn relating to (i) data; (ii) statistical models and model-fitting; and (iii) visualisation and interpretation. However, we also emphasise that these are very closely interlinked and although we have used these coarse 'pigeon-holes', there are many overlapping aspects and challenges.

### 2.1 | Data

Ecology, like environmental sciences and other branches of biology, has entered into an era of big data, with enormous possibilities for a better understanding of environmental state (Runting et al., 2022). Data can be 'big' due to different characteristics. The 'Four Vs Framework' (see discussion in Farley et al. (2018) and references therein) discuss four distinct aspects: (1) volume: quantity of data (2) velocity: time-varying data; (3) variety: multiple data types with complex relationships; and (4) veracity: trustworthiness of the data. These different aspects often do not occur in isolation, leading to multiple intricate data challenges when analysing ecological data. We highlight just some of the problems and approaches to address specific associated 'V' challenges that the authors of the papers within this Special Feature have encountered and discussed.

Biologging sensor technologies have been at the forefront of creating large volumes of available data, frequently at a range of different scales. Thus, the analysis of biologging data is often pioneering within ecology in relation to big data, with the potential to rapidly transform our understanding of the ecology, particularly in their application to animal movements (Nathan et al., 2022; Williams et al., 2020). A key limitation of most current systems, however, is the trade-off between collecting ultra-fine sub-second scale movement and behaviour data over shorter periods of time vs. more coarse but longer-term movement and space use data. Wild et al. (2023) take advantage of rapid developments in the field of the Internet of Things (i.e. methods for attaching electronic sensor devices, connected to a network, to everyday objects) to overcome key limitations in current biologging data networking systems and present new Wi-Fi solutions, combined with smart embedded software, for big biologging data. The authors are able to demonstrate orders of magnitude of improvement in data retrieval efficiency, which is the biggest limitation of animal biologging systems. In particular, Wild et al. (2023) discuss in detail challenges and solutions concerning software architecture, on-board processing of biologging sensor data, difficulties of time synchronisation and the data transmission concept and the pros and cons of different Wi-Fi infrastructures.

Advances in technology has also led to (perhaps less foreseen) forms of data gathering mechanisms gaining momentum, and associated build-up of large quantities of data, with the rise of citizen (or community) science initiatives. The resulting data from such initiatives are typically very varied in nature, often involving multiple data collection protocols with more limited/reduced structure than compared to traditional survey methods, including data arising from opportunistic events. While analysing citizen science data from designed surveys requires carefully developed methods, difficulties increase markedly with data from semi-structured projects, for example without fixed data collection protocols or data collected by observers of any degree of observer knowledge. This leads to new challenges across the whole spectrum of the 4 'V's. While these challenges have some commonality in terms of similar issues to address and overcome, due to the large expanse of types of data collection techniques, the specific challenges and associated data analytic approaches will vary. Johnston et al. (2023) summarise four overarching categories of challenges: (i) observer behaviour, including, for example spatial bias, observer or reporting differences, and false-positive errors; (ii) data structures, relating to both measures of detectability and procedures for validation; (iii) statistical models, including not only the opportunities provided by data integration and multispecies models but also sources of bias and computational limitations; and (iv) communication, motivated by the application of citizen science within biodiversity monitoring.

The veracity of data within biodiversity also arises in less obvious ways, outside the sphere of data collection protocols 'in the field', which are most commonly considered as the reason for querying the trustworthiness of the data. In particular, there is a wealth

of information contained with many ecological and biodiversity databases. However, to combine this information, data must typically be uniquely associated with specific species and taxa. This in itself raises methodological challenges, due to, for example dynamic species names, the discovery of new species, changing biological attributes, etc. As a result, homonyms, synonyms and errors may accumulate while for many taxa a general consensus on an accepted name and taxonomic and phylogenetic relationships may not have been reached so that taxonomy itself may resemble a confusingly intricate tangled bank. To address such issues, Grenié et al. (2023) provide an extensive review of the tools, databases and best practices for harmonising taxon names in biodiversity studies. In particular, they categorise the 'wild world' of existing publicly available taxonomic databases and resources, along the axes of taxonomic breadth and spatial scope, and discuss the associated strengths and caveats of each database. In addition, on the practical computation side, they review the existing computational tools provided in different R packages for taxonomic harmonisation, and, perhaps rather fittingly, provide a 'taxonomy' of the R packages, classifying them according to their associated functions.

## 2.2 | Models and model fitting

A vast array of different statistical models have been developed and fitted to ecological data in the last decade or so (Guisan et al., 2017; Hooten et al., 2017; Kery & Royle, 2016; MacKenzie et al., 2018; McCrea & Morgan, 2015; Royle et al., 2014; Schaub & Kéry, 2021), often with limited critical review of the characteristics and associated disadvantages and challenges of each. The advancement in models and associated model-fitting tools reflect the changing quantity of the data (as highlighted above), quality of the data (e.g. increased spatial/temporal resolution), emerging forms of data from new technologies (e.g. earth observation and/or drone data, eDNA) and advanced computational techniques (and associated computational power). Thus, summary overviews of these emerging and advancing areas are important and timely for ecologists and statisticians to be able to understand what can, and often importantly, what cannot (or should not), be done and also provide tools for fitting such models to different data. These models encompass all areas of ecology from population and community ecology to landscape and ecosystem ecology. Interrogation of the associated modelling ideas motivates further advances in addressing the challenges and model development to account for additional data complexities or efficient model-fitting tools, for example. We briefly summarise here some of the types of models and associated challenges that arise across a range of different types of models, and data, within this Special Feature.

Developing or adapting general statistical models that can be applied to different forms of data can be very scientifically efficient. Such approaches also often permit the use of readily available software packages, for example NIMBLE (de Valpine et al., 2017), R-INLA Lindgren and Rue (2015) and inlabru (Bachl et al., 2019) as well as

specific application focused packages, such as MARK/RMARK (for capture-recapture models; Laake, 2013); momentuHMM (for hidden Markov models [HMMs] applied to movement data; McClintock & Michelot, 2018) and Distance (for distance sampling; Thomas et al., 2010). Areas which have accessible software are witnessing substantial statistical development, enhanced by the flexibility of the computational tools provided. For example, R-INLA and inlabru have been used by both Laxton et al. (2023) and Torney et al. (2023), while Newman et al. (2023) discusses the relative merits of available software tools for fitting models. However, Barros et al. (2023) take one step further from the issue of readily accessible computer packages, suggesting that model fitting is not the primary challenge, rather that the models being used by ecologists need to be considered as predictive models, which can be used transparently and easily adapted following updated datasets or statistical methodology. Their proposal of the PERFICT workflow provides a framework by which these important challenges can be aligned.

Understanding the relationship between such general statistical models and specific ecological models can be challenging, as can be structuring the data into the required general form. Two particular 'umbrella' models that have been applied extensively within ecological models are the closely related HMMs and state-space models (SSMs). Both these types of models are widely used in ecological settings in the presence of longitudinal data (Auger-Methe et al., 2021; McClintock et al., 2021). One attraction of these models within the ecological applications is that they both directly separate out the distinct ecological and/or sampling processes. This often simplifies the model specification, permitting the consideration of the separate components independently. A common distinction between these models relates to whether the latent processes are defined to be discrete-valued (for HMMs) or continuous-valued (SSMs), although we note that this distinction is not universally used. Specific ecological areas where these models have been extensively applied, include, but are far from limited to, fisheries stock assessment (Aeberhard et al., 2018); population dynamics (Newman et al., 2014); animal movement (Hooten et al., 2017; Langrock et al., 2012; Patterson et al., 2017); and capture-recapture-type surveys (King, 2014; McCrea & Morgan, 2015). Glennie et al. (2023) and Newman et al. (2023) provide a methodological (and practical) review of HMMs and SSMs, respectively.

In particular, Glennie et al. (2023) highlight the potential difficulties that may be encountered when specifying HMMs for different systems, including issues which arise when model assumptions are not valid and the challenges of defining and fitting a suitable model in an HMM framework when the underlying hidden process increases in complexity. Providing descriptions of these general statistical models that can be applied to a variety of different forms of ecological data and associated discussion of issues to be aware of are a very useful resource for practitioners, particularly when describing the pitfalls that may arise. The rapid growth of the application of HMMs has also been aided by associated efficient model-fitting algorithms, due to the Markovian structure of the model (Zucchini et al., 2016).

The practical issues of fitting general and flexible SSMs, assuming a continuous-valued ecological (latent) process, is highlighted and addressed by Newman et al. (2023). Importantly, they discuss and contrast a wide-range of model-fitting techniques, dependent on the underlying assumptions of the specified model. In particular, they describe model-fitting algorithms that can accommodate more complex modelling dynamics, such as nonlinear processes and/or non-Gaussian stochasticity. Such models are less familiar/used within the ecological community, most likely due to the associated model-fitting challenges, however such adaptations of SSMs have great potential for the modelling of ecological data. The important aspect of what software can be used to fit such complex models is also highlighted in the paper.

The challenges of fitting models to data can concern both the associated algorithms required (as for SSMs) and the increase in computational expense, particularly as the complexity of the model increases. With increasingly large datasets, such as those routinely collected in bioacoustics or biologging studies (see Wild et al., 2023), many standard methods break down and cannot be practically applied. There is hence a necessity to identify and develop suitable modifications to improve computational efficiency and scalability, adapting traditional (and developing new) methods to big data. Providing successful examples, and the associated strategies that were most successful, including for example, computational efficiencies (Newman et al., 2023) and as demonstrated in King et al. (2022), as well as model simplifications that retain the signal within the data, are promising avenues going forward. The challenges that arise regarding scalability due to large (and new) datasets are also an opportunity for the development and use of machine learning algorithms. However, off-the-shelf algorithms may not be sufficient or may be too limiting, as described by Wang et al. (2023), so additional developments may be required for ecological applications. For example, it will generally be important to incorporate known ecological processes within the data analysis.

There are numerous opportunities, risks and trade-offs in building structurally complex models to increase insight on the underlying ecological processes. For example, Laxton et al. (2023) use the very popular species distribution models (SDMs) to highlight the importance of increasing model complexity based on ecological theory. The authors showcase the usefulness of a marked point process approach, which permits the inclusion of key population dynamic processes linked to ecological covariates (relating to landscape structure and the range of movements of the study species), and highlight the importance of maintaining an understanding of the roles and effects of each model component, to ensure interpretability and useful ecological insight. Alternatively, Torney et al. (2023) show that, in relation to the study of movement behaviour, including complex mechanisms driving animal distributions into the statistical models can substantially increase model performance and predictive ability. Furthermore, they demonstrate that the relationship between model complexity and model performance is non-monotonic, highlighting the importance of robust procedures for checking models.

## 2.3 | Interpretability and visualisation

It is now possible to fit a wealth of complex models to datasets, but where is the line drawn between fitting a model for complexity's sake and because the output is required for an understanding of the dynamics exhibited by the data? In many cases, could a simple model actually be more useful/informative? Such questions are long-standing in many areas, including ecology (Murtaugh, 2007). Statistical models continue to be developed to represent the underlying data generating ecological processes—but these will always be a simplification of reality—with more complex models aiming to extract meaningful and useful interpretable ecological insight. In general, there is a trade-off between the complexity of the model being fitted and the associated intricacy of the information that can be extracted (given suitable and available data). Furthermore, statistical learning (or machine learning) techniques are rapidly increasing in their prominence and usage within ecology (Ho & Goethals, 2022; Pichler & Hartig, 2022), with such techniques often demonstrating good predictive performance, but at the lack of ecologically interpretable parameters. It is becoming increasingly important to extract interpretable and meaningful results/output from appropriate models fitted to real data, combined with intelligent visualisations, within and beyond the wider scientific community, for example, with policy-makers

One particular area of ecology in which increasing model complexity leads to further interpretability challenges is that of species' distribution modelling. Traditionally, such models have been used to establish a correlation between a single species and the environment that it occupies in order to gain an understanding of habitat suitability, or to predict the impacts of environmental change. However, there has been growing interest for these models to go beyond a single species in isolation and to include interactions between species (Kissling et al., 2012; Pollock et al., 2014) and/or the underlying mechanisms (Buckley et al., 2010) in order to improve predictability of multispecies models. However, in increasing the complexity of the model, the associated interpretability of the model parameters can become more difficult. To address this issue, Powell-Romero et al. (2023) use a feature-based approach to describe community structure within ensemble modelling approaches to improve the practical interpretability of multispecies models. Through the inclusion of simple features to describe communities, it is possible to obtain insight of not only which models outperform others, but also why this is the case. Furthermore, within more complex dynamic SDMs, Laxton et al. (2023) argue that any increased complexity in the model needs to be grounded in ecological theory. This in turn permits greater interpretability since the different mechanisms or patterns of each component of the model can be identified leading to increased interpretable ecological insight.

As models and data become more complex and high dimensional, obtaining meaningful and useful *visualisations* of the data and/or model outputs for improved insight also becomes more challenging. Traditional methods, such as dimension reduction

and considering pair-wise correlations, may lead more nuanced and/or intricate ecological insights being masked, or even lead to biases in their presentation (McInerney et al., 2014; McInerney & Krzywinski, 2015). This is particularly challenging in more complex data/model structures, such as networks or graphs structures. For example, food web visualisation should allow us to gain an understanding of the structure of foodwebs and provide insight into the detail of the complexity; however, current approaches tend to simplify the structure and therefore cannot provide the insight needed. To address some of these challenges, Pawluczuk and Iskrzyński (2023) propose methods for visualising increasingly complex foodweb (and other network) structures by combining heatmaps, interactive and animated graphs. Alternatively, Van Moorter et al. (2023) have developed the package ConScape (in Julia) which allows users to efficiently analyse and visualise landscape and habitat connectivity more simply. Further issues arise when attempting to analyse objects that contain multiple distinct (non-independent) parts that make up the complete object (e.g. when analysing skeletons rather than individual bones). With this focus, Thomas et al. (2023) propose a method based on regularised consensus principal components analysis to be able to summarise and compare shape variation in multipart morphospaces. Importantly, they also provide an accompanying R package, to permit wider usage and impact within the large scientific community.

### 3 | CONCLUDING COMMENTS AND FUTURE OUTLOOK

The opportunities for gaining an understanding of ecological systems from the range of different forms of available data (and new emerging data) are immense. However, to fully capitalise on these opportunities, addressing the associated challenges and achieving academic and societal impact, a multidisciplinary approach considering the whole data analytic pipeline is required. We discuss a number of important aspects that will contribute to advancing ecological knowledge and address important societal issues (though we note that this is far from an exhaustive list):

#### 3.1 | Interdisciplinarity

Immersive interdisciplinarity in the ecological community's research approach has the largest potential for achieving research step-changes within the discipline. The cross-fertilisation of knowledge from, for example ecologists, engineers (designing data collection devices), statisticians (developing advanced modelling techniques to fully exploit the available data and designing survey sampling strategies) and computer scientists (offering expertise in machine learning and automation) provides the opportunity for the co-creation of new and exciting approaches to address challenging ecological problems. Close collaboration with mathematical ecologists allows a better realistic connection of models

to ecological theory; equally important is the collaboration with ecologists at the model output stage to build confidence that the results are biologically realistic.

#### 3.2 | Data-centric methodological innovation

It is important to ensure that data analytic methods are being developed to make the most of the diverse and sizeable amounts of ecological data now being efficiently collected at increasing scale and quantity (Zipkin et al., 2021). However, the advancement of data collection technology continues at a rapid pace, and necessarily the associated data analytic tools are developed at a lagged timescale (there is no point in developing analytic tools for data that do not exist and/or cannot be collected). Again, an interdisciplinary outlook will help identifying novel data collection tools and methods not used yet in ecology.

#### 3.3 | Robust data integration

There has been a natural development towards integrating datasets within a single model in recent years (Frost et al., 2023), spanning both multiple data types of a single species (Isaac et al., 2020) and data from multiple species (Barraquand & Gimenez, 2019). This means that one of the biggest challenges facing statistical ecologists is to think about whether the types of data being combined in an analysis are indeed comparable—do they have differing quality and will this affect the model performance? For example, will combining small structured datasets with large unstructured data, for example from the Global Biodiversity Information Facility (GBIF), help to limit the bias in the latter, or the context dependency in the former (Isaac et al., 2020)?

#### 3.4 | 'All models are wrong, but some are useful'

This phrase, attributed to the statistician George Box, continues to provide useful insight. In particular, we apply this reasoning to the idea that being able to fit complex statistical models to data (accessible through advances in associated software) does not mean that the models are appropriate (or useful) for the data. There is a need to consider the philosophy of 'should we' fit a model to a given dataset, and ask whether it is necessary and/or appropriate given the particular ecological question of interest and available data. Gain in knowledge should trump model complexity or methods sophistication per se.

#### 3.5 | Machine learning and artificial intelligence

Such approaches are likely to have an important role in the future direction of methods in the ecological domain (Pichler & Hartig, 2022),



particularly when prediction is a primary objective. However, such methods should not simply be blindly applied to align with popular analytical trends—it is important that there is a methodological driver underpinning their usage. The interpretability of such models is more challenging due to the ‘black-box’ nature of the algorithms and lack of ecological constraints or input, for example. Considerable debate and uncertainty remains in the validity and best practices of these approaches particularly in relation to generalisability, conceptual simplicity, robustness and transparency. There is a need to increase research efforts into machine learning and artificial intelligence approaches so that their power can be appropriately harnessed for ecology and evolution. For example, novel understanding from carefully fitted and interpreted machine learning methods could be more often also used to guide the development of new likelihood-based methods.

### 3.6 | Software

This is an increasingly prominent feature of statistical analyses. The type of software ranges from general statistical packages to which ecological models and data analyses can be conducted (such as Inlabru Bachl et al., 2019 or NIMBLE de Valpine et al., 2017), to specialised packages for very specific problems (Van Moorter et al., 2023). However, the variety of computer packages (and in different languages, such as R or Python or Julia) leads to additional challenges of identifying the most relevant and/or efficient for the given problem at hand. Clear guidance regarding the advantages and disadvantages of different approaches is a particularly useful resource, though often difficult as there may be many different data and question dependent decisions in practice.

### 3.7 | Communication

The importance of improved communication for addressing and solving the inherent challenges of citizen science data are highlighted in Johnston et al. (2023). In particular, the authors focus on the importance of disseminating new statistical methods beyond the limited circle of technical groups. This requires moving beyond code sharing, investing also in software development and teaching activities and resources. They also conclude that a ‘democratisation’ of data analysis may emulate the progress brought by the democratisation of data collection through citizen science and help make the most of these data, which has to be one of the most pressing issues facing statistical ecologists at this current time.

The papers in this Special Feature only scratch the surface of the challenges present with large data and complex models, and propose some possible approaches for dealing with different issues and advance our ecological understanding. These areas of research will continue to provide a rich and diverse set of challenges for ecological researchers, but recognising the challenges, building interdisciplinary data analytic pipelines and providing interpretable results will ensure the research produced by this cross-disciplinary academic

community will reach its full potential, leading to step-changes in our ecological understanding, and be a firm basis for informed policy decision-making.

### ACKNOWLEDGEMENTS

This special feature arose from discussions and interactions at the 2019 ICMS-funded meeting ‘Addressing Statistical Challenges of Modern Technological Advances’, organised by the National Centre for Statistical Ecology, and the joint BES Quantitative and Movement Ecology Special Interest Group Meeting in Sheffield in 2018. RM is currently funded by EPSRC grant EP/S020470/1.

### PEER REVIEW

The peer review history for this article is available at <https://www.webofscience.com/api/gateway/wos/peer-review/10.1111/2041-210X.14050>.

### ORCID

Rachel McCrea  <https://orcid.org/0000-0002-3813-5328>

Ruth King  <https://orcid.org/0000-0002-5174-8727>

Laura Graham  <https://orcid.org/0000-0002-3611-7281>

Luca Börger  <https://orcid.org/0000-0001-8763-5997>

### REFERENCES

- Aeberhard, W. H., Flemming, J. M., & Nielsen, A. (2018). Review of state-space models for fisheries science. *Annual Review of Statistics and Its Application*, 5, 215–235.
- Auger-Methe, M., Newman, K., Cole, D., Empacher, F., Gryba, R., King, A. A., Leos-Barajas, V., Flemming, J. M., Nielsen, A., Petris, G., & Thomas, L. (2021). A guide to state–space modeling of ecological time series. *Ecological Monographs*, 91, 1–38.
- Bachl, F. E., Lindgren, F., Borchers, D. L., & Illian, J. B. (2019). Inlabru: An R package for Bayesian spatial modelling from ecological survey data. *Methods in Ecology and Evolution*, 10(6), 760–766.
- Barraquand, F., & Gimenez, O. (2019). Integrating multiple data sources to fit matrix population models for interacting species. *Ecological Modelling*, 411, 108713.
- Barros, C., Luo, Y., Chubaty, A., Eddy, I., Micheletti, T., Boisvenue, C., Andison, D., Cumming, S., & McIntire, E. (2023). Empowering ecological modellers with a PERFICT work-flow: Seamlessly linking data, parameterisation, prediction, validation and visualisation. *Methods in Ecology and Evolution*, 14(1), 173–188.
- Buckley, L. B., Urban, M. C., Angilletta, M. J., Crozier, L. G., Rissler, L. J., & Sears, M. W. (2010). Can mechanism inform species' distribution models? *Ecology Letters*, 13(8), 1041–1054.
- de Valpine, P., Turek, D., Paciorek, C. J., Anderson-Bergman, C., Lang, D. T., & Bodik, R. (2017). Programming with models: Writing statistical algorithms for general model structures with NIMBLE. *Journal of Computational and Graphical Statistics*, 26(2), 403–413.
- Farley, S. S., Dawson, A., Goring, S. J., & Williams, J. W. (2018). Situating ecology as a big-data science: Current advances, challenges, and solutions. *Bioscience*, 68(8), 563–576.
- Fer, I., Kelly, R., Moorcroft, P. R., Richardson, A. D., Cowdery, E. M., & Dietze, M. C. (2018). Linking big models to big data: Efficient ecosystem model calibration through Bayesian model emulation. *Biogeosciences*, 15(19), 5801–5830.
- Frost, F., McCrea, R. S., King, R., Gimenez, O., & Zipkin, E. (2023). Integrated population models: Achieving their potential. *Journal of Statistical Theory and Practice*, 17, 6.

- Glennie, R., Adam, T., Leos-Barajas, V., Michelot, T., Photopoulou, T., & McClintock, B. T. (2023). Hidden markov models: Pitfalls and opportunities in ecology. *Methods in Ecology and Evolution*, 14(1), 43–56.
- Grenié, M., Berti, E., Carvajal-Quintero, J., Dädlow, G. M. L., Sagouis, A., & Winter, M. (2023). Harmonizing taxon names in biodiversity data: A review of tools, databases and best practices. *Methods in Ecology and Evolution*, 14(1), 12–25.
- Guisan, A., Thuiller, W., & Zimmermann, N. E. (2017). *Habitat suitability and distribution models: With applications in R*. Cambridge University Press.
- Ho, L., & Goethals, P. (2022). Machine learning applications in river research: Trends, opportunities and challenges. *Methods in Ecology and Evolution*, 13(11), 2603–2621.
- Hooten, M. B., Johnson, D. S., McClintock, B. T., & Morales, J. M. (2017). *Animal movement: Statistical models for telemetry data*. CRC Press.
- Isaac, N. J. B., Jarzyna, M. A., Keil, P., Dambly, L. I., Boersch-Supan, P. H., Browning, E., Freeman, S. N., Golding, N., Guillera-Aroita, G., Henrys, P. A., Jarvis, S., Lahoz-Monfort, J., Pagel, J., Pescott, O. L., Schmucki, R., Simmonds, E. G., & O'Hara, R. B. (2020). Data integration for large-scale models of species distributions. *Trends in Ecology & Evolution*, 35(1), 56–67.
- Johnston, A., Matechou, E., & Dennis, E. B. (2023). Outstanding challenges and future directions for biodiversity monitoring using citizen science data. *Methods in Ecology and Evolution*, 14(1), 103–116.
- Kery, M., & Royle, J. A. (Eds.). (2016). *Applied hierarchical modeling in ecology*. Academic Press.
- King, R. (2014). Statistical ecology. *Annual Review of Statistics and its Application*, 1, 410–426.
- King, R., Sarzo, B., & Elvira, V. (2022). *When ecological individual heterogeneity models and large data collide: An importance sampling approach* (Technical report). University of Edinburgh. <https://arxiv.org/abs/2205.07261>
- Kissling, W. D., Dormann, C. F., Groeneveld, J., Hickler, T., Kühn, I., McInerney, G. J., Montoya, J. M., Römermann, C., Schiffers, K., Schurr, F. M., Singer, A., Svenning, J.-C., Zimmermann, N. E., & O'Hara, R. B. (2012). Towards novel approaches to modelling biotic interactions in multispecies assemblages at large spatial extents. *Journal of Biogeography*, 39(12), 2163–2178.
- Laake, J. (2013). RMark: An R interface for analysis of capture-recapture data with MARK. AFSC processed rep. 2013-01, Alaska Fisheries Science Center, NOAA, National Marine Fisheries Service, Seattle, WA.
- Langrock, R., King, R., Matthiopoulos, J., Thomas, L., Fortin, D., & Morales, J. M. (2012). Flexible hidden Markov-type models for animal telemetry data. *Ecology*, 93, 2336–2342.
- Laxton, M. R., de Rivera, O. R., Soriano-Redondo, A., & Illian, J. B. (2023). Balancing structural complexity with ecological insight in spatio-temporal species distribution models. *Methods in Ecology and Evolution*, 14(1), 162–172.
- Lindgren, F., & Rue, H. (2015). Bayesian spatial modelling with R-INLA. *Journal of Statistical Software*, 63(19), 1–25.
- MacKenzie, D. I., Nichols, J. D., Royle, J. A., Pollock, K. H., Bailey, L. L., & Hines, J. E. (2018). *Occupancy estimation and modeling* (2nd ed.). Academic Press.
- McClintock, B. T., Langrock, R., Gimenez, O., Cam, E., Borchers, D. L., Glennie, R., & Patterson, T. A. (2021). Uncovering ecological state dynamics with hidden Markov models. *Ecology Letters*, 23, 1878–1903.
- McClintock, B. T., & Michelot, T. (2018). momentuHMM: R package for generalized hidden Markov models of animal movement. *Methods in Ecology and Evolution*, 9(6), 1518–1530.
- McCrea, R. S., & Morgan, B. J. T. (2015). *Analysis of capture-recapture data*. Chapman and Hall/CRC Press.
- McInerney, G., & Krzywinski, M. (2015). Unentangling complex plots. *Nature Methods*, 12(7), 591. Number: 7 Publisher: Nature Publishing Group.
- McInerney, G. J., Chen, M., Freeman, R., Gavaghan, D., Meyer, M., Rowland, F., Spiegel-Halter, D. J., Stefaner, M., Tassarolo, G., & Hortal, J. (2014). Information visualisation for science and policy: Engaging users and avoiding bias. *Trends in Ecology & Evolution*, 29(3), 148–157.
- Murtaugh, P. A. (2007). Simplicity and complexity in ecological data analysis. *Ecology*, 88(1), 56–62.
- Nathan, R., Monk, C. T., Arlinghaus, R., Adam, T., Alós, J., Assaf, M., Baktoft, H., Beardsworth, C. E., Bertram, M. G., Bijleveld, A. I., Brodin, T., Brooks, J. L., Campos-Candela, A., Cooke, S. J., Gjelland, K., Gupte, P. R., Harel, R., Hellström, G., Jeltsch, F., ... Jarić, I. (2022). Big-data approaches lead to an increased understanding of the ecology of animal movement. *Science*, 375(6582), eabg1780.
- Newman, K., King, R., Elvira, V., de Valpine, P., McCrea, R. S., & Morgan, B. J. T. (2023). State-space models for ecological time-series data: Practical model-fitting. *Methods in Ecology and Evolution*, 14(1), 26–42.
- Newman, K. B., Buckland, S. T., Morgan, B. J. T., King, R., Borchers, D. L., Cole, D. J., Besbeas, P., Gimenez, O., & Thomas, L. (2014). *Modelling population dynamics: Model formulation, fitting and assessment using state-space methods*. Springer.
- Niu, S., Wang, S., Wang, J., Xia, J., & Yu, G. (2020). Integrative ecology in the era of big data—From observation to prediction. *Science China Earth Sciences*, 63, 1429–1442.
- Patterson, T. A., Parton, A., Langrock, R., Blackwell, P. G., Thomas, L., & King, R. (2017). Statistical modelling of individual animal movement: An overview of key methods and a discussion of practical challenges. *Advances in Statistical Analysis*, 101, 399–438.
- Pawluczuk, Ł., & Iskrzyński, M. (2023). Food web visualisation: Heat map, interactive graph and animated flow network. *Methods in Ecology and Evolution*, 14(1), 57–64.
- Pichler, M., & Hartig, G. (2022). Machine learning and deep learning – A review for ecologists. Technical report. <https://arxiv.org/abs/2204.05023>
- Pollock, L. J., Tingley, R., Morris, W. K., Golding, N., O'Hara, R. B., Parris, K. M., Vesk, P. A., & McCarthy, M. A. (2014). Understanding co-occurrence by modelling species simultaneously with a joint species distribution model (JSDM). *Methods in Ecology and Evolution*, 5(5), 397–406.
- Powell-Romero, F., Fountain-Jones, N. M., Norberg, A., & Clark, N. J. (2023). Improving the predictability and interpretability of co-occurrence modelling through feature-based joint species distribution ensembles. *Methods in Ecology and Evolution*, 14(1), 146–161.
- Reyers, B., & Selig, E. R. (2020). Global targets that reveal the social-ecological interdependencies of sustainable development. *Nature Ecology and Evolution*, 4, 1011–1019.
- Royle, J., Chandler, R. B., Sollmann, R., & Gardner, B. (2014). *Spatial Capture-recapture*. Academic Press.
- Runting, R. K., Phinn, S., Xie, Z., Veter, O., & Watson, J. E. M. (2022). Opportunities for big data in conservation and sustainability. *Nature Communications*, 11, 2003.
- Schaub, M., & Kéry, M. (2021). *Integrated population models*. Academic Press.
- Solé, R., & Levin, S. (2022). Ecological complexity and the biosphere: The next 30 years. *Philosophical Transactions of the Royal Society B: Biological Sciences*, 377(1857), 20210376.
- Sutherland, W. J., Freckleton, R. P., Godfray, H. C. J., Beissinger, S. R., Benton, T., Cameron, D. D., Carmel, Y., Coomes, D. A., Coulson, T., Emmerson, M. C., Hails, R. S., Hays, G. C., Hodgson, D. J., Hutchings, M. J., Johnson, D., Jones, J. P. G., Keeling, M. J., Kokko, H., Kunin, W. E., ... Wiegand, T. (2013). Identification of 100 fundamental ecological questions. *Journal of Ecology*, 101(1), 58–67.
- Thomas, D. B., Harmer, A. M. T., Giovanardi, S., Holvast, E. J., McGovern, C. M., & Tenenhaus, A. (2023). Constructing a multiple-part morphospace using a multiblock method. *Methods in Ecology and Evolution*, 14(1), 65–76.



- Thomas, L., Buckland, S. T., Rexstad, E. A., Laake, J. L., Strindberg, S., Hedley, J. S. L., Bishop, R. B., Marques, T. A., & Burnham, K. P. (2010). Distance software: Design and analysis of distance sampling surveys for estimating population size. *Journal of Applied Ecology*, *47*, 5–14.
- Torney, C. J., Laxton, M., Lloyd-Jones, D. J., Kohi, E. M., Frederick, H. L., Moyer, D. C., Mrisha, C., Mwita, M., & Hopcraft, J. G. C. (2023). Estimating the abundance of a group-living species using multi-latent spatial models. *Methods in Ecology and Evolution*, *14*(1), 77–86.
- Van Moorter, B., Kivimäki, I., Noack, A., Devooght, R., Panzacchi, M., Hall, K. R., Leleux, P., & Saerens, M. (2023). Accelerating advances in landscape connectivity modelling with the ConScape library. *Methods in Ecology and Evolution*, *14*(1), 133–145.
- Wang, Z., Gong, H., Huang, M., Gu, F., Wei, J., Guo, Q., & Song, W. (2023). A multi-model random forest ensemble method for an improved assessment of chinese terrestrial vegetation carbon density. *Methods in Ecology and Evolution*, *14*(1), 117–132.
- Wild, T. A., Wikelski, M., Tyndel, S., Alarcón-Nieto, G., Klump, B. C., Aplin, L. M., Meboldt, M., & Williams, H. J. (2023). Internet on animals: Wi-fi-enabled devices provide a solution for big data transmission in biologging. *Methods in Ecology and Evolution*, *14*(1), 87–102.
- Williams, H. J., Taylor, L. A., Benhamou, S., Bijleveld, A. I., Clay, T. A., de Grissac, S., Demšar, U., English, H. M., Franconi, N., Gómez-Laich, A., Griffiths, R. C., Kay, W. P., Morales, J. M., Potts, J. R., Rogerson, K. F., Rutz, C., Spelt, A., Trevail, A. M., Wilson, R. P., & Börger, L. (2020). Optimizing the use of biologgers for movement ecology research. *Journal of Animal Ecology*, *89*(1), 186–206.
- Yates, K. L., Bouchet, P. J., Caley, M. J., Mengersen, K., Randin, C. F., Parnell, S., Fielding, A. H., Bamford, A. J., Ban, S., Barbosa, A. M., Dormann, C. F., Elith, J., Embling, C. B., Ervin, G. N., Fisher, R., Gould, S., Graf, R. F., Gregr, E. J., Halpin, P. N., ... Sequeira, A. M. (2018). Outstanding challenges in the transferability of ecological models. *Trends in Ecology Evolution*, *33*(10), 790–802.
- Zipkin, E. F., Zylstra, E. R., Wright, A. D., Saunders, S. P., Finley, A. O., Dietze, M. C., Itter, M. S., & Tingley, M. W. (2021). Addressing data integration challenges to link ecological processes across scales. *Frontiers in Ecology and the Environment*, *19*(1), 30–38.
- Zucchini, W., MacDonald, I., & Langrock, R. (2016). *Hidden Markov models for time series: An introduction using R* (2nd ed.). Chapman and Hall/CRC.