

Kent Academic Repository

Full text document (pdf)

Citation for published version

Diana, Alex, Matechou, Eleni, Griffin, Jim E., Buxton, Andrew S. and Griffiths, Richard A. (2021) An RShiny app for modelling environmental DNA data: accounting for false positive and false negative observation error. *Ecography*, 44 . pp. 1838-1844. ISSN 0906-7590.

DOI

<https://doi.org/10.1111/ecog.05718>

Link to record in KAR

<https://kar.kent.ac.uk/91832/>

Document Version

Publisher pdf

Copyright & reuse

Content in the Kent Academic Repository is made available for research purposes. Unless otherwise stated all content is protected by copyright and in the absence of an open licence (eg Creative Commons), permissions for further reuse of content should be sought from the publisher, author or other copyright holder.

Versions of research

The version in the Kent Academic Repository may differ from the final published version.

Users are advised to check <http://kar.kent.ac.uk> for the status of the paper. **Users should always cite the published version of record.**

Enquiries

For any further enquiries regarding the licence status of this document, please contact:

researchsupport@kent.ac.uk

If you believe this document infringes copyright then please contact the KAR admin team with the take-down information provided at <http://kar.kent.ac.uk/contact.html>

Software notes

An RShiny app for modelling environmental DNA data: accounting for false positive and false negative observation error

Alex Diana, Eleni Matechou, Jim E. Griffin, Andrew S. Buxton and Richard A. Griffiths

A. Diana and E. Matechou (<https://orcid.org/0000-0003-3626-844X>) ✉ (e.matechou@kent.ac.uk), School of Mathematics, Statistics and Actuarial Science, Univ. of Kent, Canterbury, UK. – J. E. Griffin, Dept of Statistical Science, Univ. College London, London, UK. – A. S. Buxton (<https://orcid.org/0000-0002-0555-2491>) and R. A. Griffiths (<https://orcid.org/0000-0002-5533-1013>), Durrell Inst. of Conservation and Ecology, School of Anthropology and Conservation, Univ. of Kent, Canterbury, UK.

Ecography

44: 1838–1844, 2021

doi: 10.1111/ecog.05718

Subject Editor: Brody Sandel
Editor-in-Chief: Miguel Araújo
Accepted 2 September 2021

Environmental DNA (eDNA) surveys have become a popular tool for assessing the distribution of species. However, it is known that false positive and false negative observation error can occur at both stages of eDNA surveys, namely the field sampling stage and laboratory analysis stage. We present an RShiny app that implements the Griffin et al. (2020) statistical method, which accounts for false positive and false negative errors in both stages of eDNA surveys that target single species using quantitative PCR methods. Following Griffin et al. (2020), we employ a Bayesian approach and perform efficient Bayesian variable selection to identify important predictors for the probability of species presence as well as the probabilities of observation error at either stage. We demonstrate the RShiny app using a data set on great crested newts collected by Natural England in 2018, and we identify water quality, pond area, fish presence, macrophyte cover and frequency of drying as important predictors for species presence at a site. The state-of-the-art statistical method that we have implemented is the only one that has specifically been developed for the purposes of modelling false negative and false positive observation error in eDNA data. Our RShiny app is user-friendly, requires no prior knowledge of R and fits the models very efficiently. Therefore, it should be part of the tool-kit of any researcher or practitioner who is collecting or analysing eDNA data.

Keywords: Bayesian variable selection, environmental DNA, multi-level occupancy model, PCR



Background

Environmental DNA (eDNA) is increasingly used within biodiversity assessments (McClenaghan et al. 2020). The method relies on the detection of DNA released from source organisms into aquatic or terrestrial environments. This DNA is extracted from a sample of the substrate, usually water or soil (Thomsen and Willerslev 2015) (stage 1), and then analysed using qPCR (Thomsen et al. 2012) or metabarcoding (Valentini et al. 2016) (stage 2).



We present an RShiny app, ‘eDNA 1.0’, for modelling single-species eDNA data resulting from qPCR analysis by implementing the Bayesian model developed by Griffin et al. (2020). The app can be used to model *eDNA scores*, which are defined as the number of qPCR runs that have successfully been amplified, for each sample and site. The model estimates site-specific probabilities of species presence while accounting for false positive and false negative observation error at both stages of eDNA surveys. The Griffin et al. (2020) model is an extension of the work by Guillera-Arroita et al. (2017) but in contrast to the latter, the Griffin et al. (2020) model does not require augmenting the eDNA data with other types of survey data. This is due to the specification of novel informative prior distributions that reflect our belief that the probability of a false positive observation is smaller than the probability of a true positive observation at either stage. Nevertheless, if opportunistic records of species presence exist for any of the sites, these can easily be accounted for within the model. Finally, Griffin et al. (2020) presented an MCMC algorithm that employs the Pólya-Gamma sampling scheme (Polson et al. 2013) and hence enables fast computation times and efficient Bayesian variable selection for all model parameters.

Methods and features

Bayesian model

Griffin et al. (2020) presented a hierarchical Bayesian model that describes the different stages of eDNA surveys in terms of the probabilities of species presence and the probabilities of observation error. All of these probabilities can be functions of site-specific covariates, with dependencies modelled using logistic regressions. Subscript s is used to denote sites, $s = 1, \dots, S$, while subscript m is used to denote samples from sites, $m = 1, \dots, M$. The list of parameters is given in Table 1 and a schematic representation of the model is provided in Fig. 1.

As explained in Griffin et al. (2020), the model is only locally identifiable, so that there exist countably many, in this case four, sets of parameter values that result in the same likelihood function value (Cole 2020). This identifiability issue is overcome by introducing informative prior distributions that express our belief that a false positive observation is less likely than a corresponding true positive at each stage, so

that the probabilities that $\theta_{11} < \theta_{10}$ or $p_{11} < p_{10}$ are small. We clarify here that these constraints apply to the probabilities of observation error and regardless of the probability of species presence. For example, we expect that θ_{11} is greater than θ_{10} , which suggests that the probability that a sample from an *occupied* site includes DNA of the targeted species is greater than the probability that a sample from an *unoccupied* site includes DNA of the targeted species, and similarly for the probabilities referring to the results of the qPCR analysis.

The MCMC algorithm presented in Griffin et al. (2020) employs the Pólya-Gamma sampling scheme (Polson et al. 2013), enabling fast computation for logistic regression models and efficient Bayesian variable selection, which is performed using an Add–Delete–Swap algorithm (Brown et al. 1998, Chipman et al. 2001). As the name of the algorithm suggests, at each MCMC iteration, we either propose to add a covariate that is not currently in the model, to delete a covariate that is in the model or to swap a covariate that is in the model with one that is not in the model at that iteration. This process gives rise to the posterior inclusion probabilities (PIPs, Barbieri et al. 2004), which indicate the proportion of iterations that each covariate was in the model for each parameter. PIPs can be used to understand how useful each of the covariates is as a predictor for the corresponding parameter and often a threshold of 0.5 is applied to identify the most important predictors (Ghosh 2015).

RShiny app

The RShiny app ‘eDNA 1.0’ is freely available and can be accessed via the RShiny server <<https://seak.shinyapps.io/eDNA/>>. However, we recommend that the app is downloaded via our dedicated website <<https://blogs.kent.ac.uk/edna/>> and run locally. Our website includes information on how to download the app, in the Download tab, as well as simulated data and a step-by-step analysis of a simulated data set, in the Examples tab. The app is implemented in R (<www.r-project.org>), with several functions written in C++ for faster implementation. It also relies on a number of shiny packages: shiny (Chang et al. 2020), shinythemes (Chang 2018), shinycssloaders (Sali and Attali 2020), shinyalert (Sali and Attali 2020) and shinyjs (Attali 2020), packages for creating plots: ggplot (Wickham 2016) and grid (<www.r-project.org>), for running and summarising posterior simulation: coda (Plummer et al. 2006), for manipulating data:

Table 1. Parameters of the Griffin et al. (2020) model. Note that $\theta_{01s} = 1 - \theta_{11s}$, $\theta_{00s} = 1 - \theta_{10s}$, $p_{01s} = 1 - p_{11s}$, $p_{00s} = 1 - p_{10s}$, and hence our RShiny app only reports results in terms of the probabilities of a positive (either true or false) observation at either stage.

	Name (probability of)	Detailed explanation (probability that)
ψ_s	Species presence	Site s is occupied by the target species
Stage 1		
θ_{11s}	Stage 1 true positive observation	A sample from occupied site s includes DNA of the target species
θ_{10s}	Stage 1 false positive observation	A sample from unoccupied site s includes DNA of the target species
Stage 2		
p_{11s}	Stage 2 true positive observation	A qPCR on a sample (from site s) that includes DNA of the target species is positive
p_{10s}	Stage 2 false positive observation	A qPCR on a sample (from site s) that does not include DNA of the target species is positive

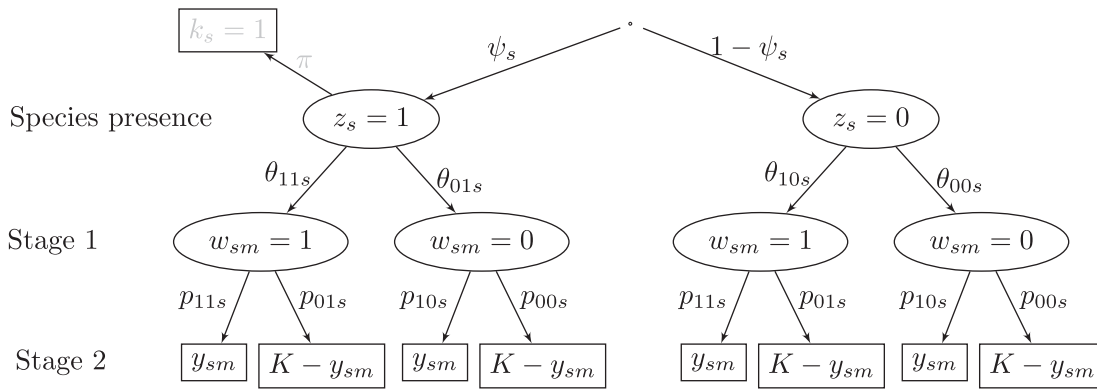


Figure 1. Schematic representation of the Griffin et al. (2020) model. Unobservable states are represented by ellipses and data by rectangles. The parameters are defined in Table 1. The latent variable z_s indicates whether site s is occupied by the target species (1) or not (0) and the latent variable w_{sm} indicates whether sample m from site s includes DNA of the target species (1) or not (0). The part of the model that is presented in grey corresponds to how opportunistic records of species presence are modelled. Specifically, parameter π indicates the probability that an occupied site has an opportunistic record associated with it and indicator variable k_s indicates whether site s is known to be occupied (1) or not (0).

reshape (Wickham 2007) and dplyr (Wickham et al. 2020) and finally, for making a beeping sound when the MCMC has finished running: beepR (Bååth 2018).

The app includes a detailed help section with several tabs that provides a step-by-step description of how to format the data, which need to be uploaded in a .csv file, how to upload the data, how to fit the model and how to access the results. Once the data have been uploaded, the user needs to specify the number of qPCR runs for each sample (this is assumed to be the same for all samples) and to select the parameters that are to be considered as functions of covariates, as shown in Fig. 2. It is not necessary to consider covariates for any of the parameters, but the set of covariates uploaded will be considered as potential predictors for all parameters that have been specified in the settings window.

The settings window also allows users to change the number of iterations, including number of chains, burn-in

and thinning, as well as the prior distribution parameters, although we would recommend that the prior settings are not changed unless the user has a good understanding of the model. Once the user clicks on Run, the app will begin model fitting and the iteration number will be shown in the bottom right corner.

The results are available in the Results tab. These include posterior summaries for all model parameters, or corresponding coefficients of covariates for parameters that have been modelled as functions of covariates. All of the results and figures that are produced as part of the output can be downloaded.

The diagnostics tab produces traceplots for all model parameters as well as effective sample sizes (ESS) obtained. A message will appear to indicate if any of the parameters have ESS lower than 500 so a closer inspection of the traceplots and ESS outputted would help identify the parameters that are not mixing well.

Figure 2. Settings tab in the eDNA RShiny app.

Table 2. List of covariates considered as predictors for the probability of occupancy of great crested newts in the example. These are based on the Habitat suitability index proposed by Oldham et al. (2000).

Covariate	Type	Levels
Geographic location	categorical	1 – optimal; 2 – marginal; 3 – unsuitable
Frequency of pond drying	categorical	1 – never; 2 – rarely; 3 – sometimes; 4 – annually
Water quality	categorical	1 – bad; 2 – poor; 3 – moderate; 4 – good
Waterfowl intensity	categorical	1 – absent; 2 – minor; 3 – major
Fish intensity	categorical	1 – absent; 2 – possible; 3 – minor; 4 – major
Terrestrial habitat quality	categorical	1 – bad; 2 – poor; 3 – moderate; 4 – good
Percentage pond shading	numerical	NA
Pond area	numerical	NA
Pond density	numerical	NA
Percentage macrophyte cover	numerical	NA

Example

We consider a data set on great crested newts collected by Natural England in 2018. $M=1$ water sample was collected from each of $S=2215$ sites and $K=12$ qPCR runs were performed for each water sample. A replicate was considered positive if an exponential growth phase was observed within the qPCR amplification curve. We have considered six categorical covariates and four continuous covariates (listed in Table 2) for the probability of occupancy, while all other parameters have been modelled as constant. We have also accounted for confirmed species presences that were available for 120 sites (see Fig. 1 for description on how opportunistic data of this type are modelled). We run 1000 burn-in iterations and 2000 additional iterations with thinning set to 20. Fitting the model using the RShiny server took just under 2 h, despite the large number of sites and considerable number of covariates considered. We note that running time depends on the number of sites, number of samples per site, number of covariates considered and of course operating system.

The app outputs posterior summaries of the site-specific probability of species presence, saved in a .csv file. For illustration purposes, we plot these summaries for a random sample of 50 sites in Fig. 3 and provide the code for producing similar plots in the Help section of the app.

The PIPs for the probability of occupancy are provided by the app in a plot (Fig. 4). A high PIP indicates stronger support for a covariate as a predictor. In this case, water quality, pond area, fish presence, percentage of macrophyte cover and frequency of pond drying all stand out as important predictors for the probability that a pond is occupied by great crested newts.

Posterior summaries of the corresponding coefficients are given in Fig. 5, showing that, even though geographic location and the percentage of pond that is shaded had PIPs above the 0.5 threshold, the 95% posterior credible intervals (PCIs) for the corresponding coefficients include 0. On the other hand, we can see that better water quality, lower pond area, lower levels of fish presence, higher levels of macrophytes and pond desiccation increase the probability of a pond being occupied by great crested newts. The predictors identified by the model and their corresponding effects are broadly consistent with current understanding of the preference of great

crested newts for vegetated, fish-free and clean water ponds (Oldham et al. 2000). These results demonstrate that important predictors for the probability of species presence can be identified using eDNA data and our RShiny app.

Posterior summaries of the probabilities related to observation error in both stages are given in Table 3. Stage 1 is related to higher probabilities of false observations, either positive or negative, compared to stage 2. The processes by which samples are collected mean that it is more likely for DNA to fail to be collected in the field, or contamination to be introduced at this stage, while lab protocols are more tightly controlled.

Finally, the app outputs the posterior probability of species absence conditional on $x=0, \dots, K$ positive qPCR replicates. For this example, the results are shown in the first row of Table 4 where we can see that the posterior conditional probability of species absence is very close to 1 given four or fewer qPCR positives, but then it declines sharply and plateaus at

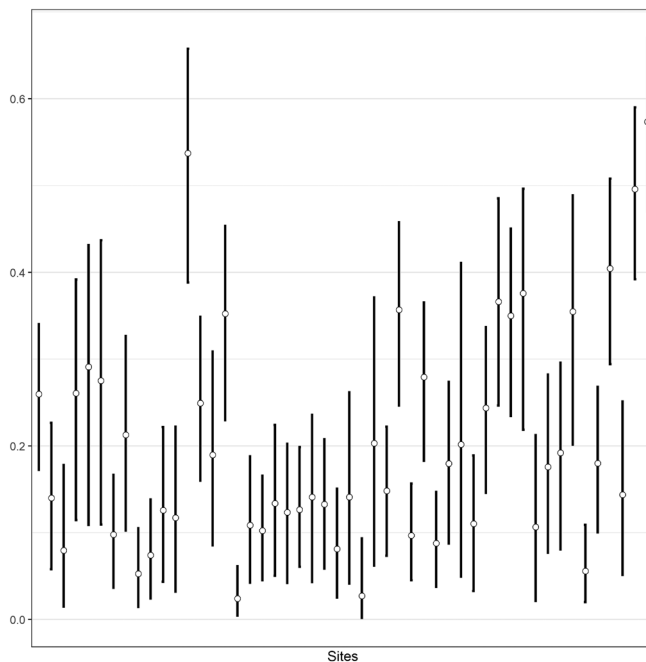


Figure 3. Posterior summaries of site-specific probabilities of occupancy for a random sample of 50 sites.

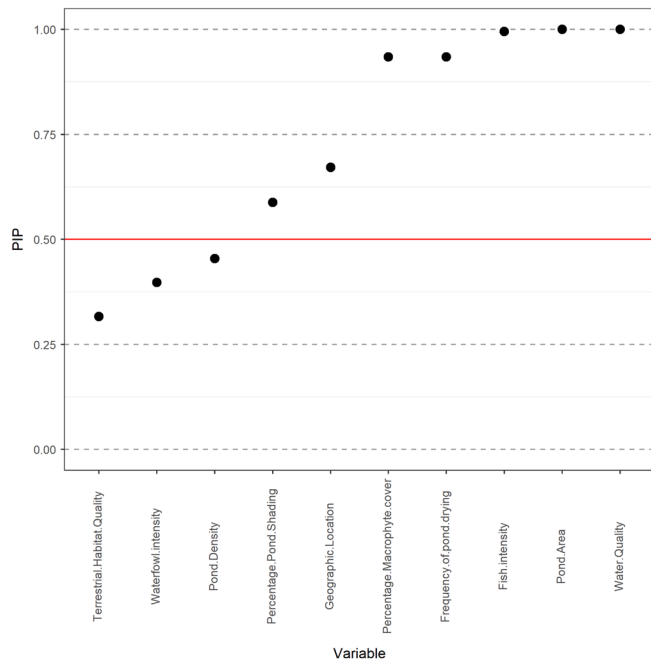


Figure 4. PIPs for the probability of occupancy. The horizontal line indicates the PIP = 0.5 line.

around 37%. The second row of Table 4 shows the posterior probability of x , $x=0, \dots, 12$, positive qPCR replicates conditional on species presence. This conditional distribution is clearly bimodal. Specifically, the posterior probability of zero qPCR positives given species presence is just under 10% and this probability decreases for $x=1, 2, 3, 4$ before it starts to increase again reaching the second peak at $x=11$ (28%). This bimodality is due to the observation error in stage 1: the first peak at 0 is a result of a stage 1 false negative observation, whereas the second peak is a result of a stage 1 true positive observation.

It is important to note that when $M=1$ the model is non-identifiable and hence the results obtained are not reliable, unless the probability of occupancy, ψ , or the probabilities of observation error in stage 1, θ_{11} and θ_{10} , are modelled as functions of covariates. Incorporating covariates helps overcome the identifiability issues, while an alternative solution is to incorporate information on confirmed species presences at some of the sites. These confirmed species presences can be, for example, opportunistic records of species presence while collecting samples from sites. Similarly, when $K=1$, the probabilities of observation error, either in stage 1 (θ_{11} and θ_{10}) or in stage 2 (p_{11} and p_{10}), need to be functions of covariates for the model to be identifiable.

Discussion

As eDNA surveys become increasingly used as monitoring tools, they have the potential to replace traditional survey methods that rely on direct observation of species, especially for difficult to detect species. Our RShiny app provides the

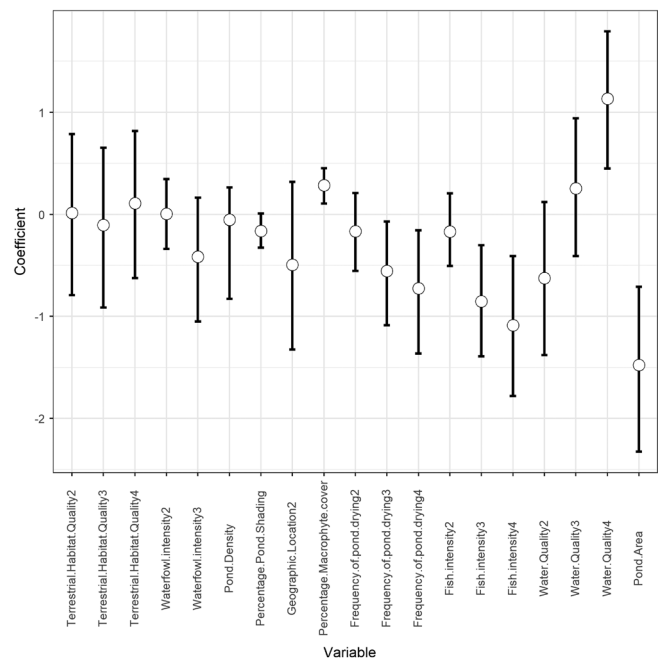


Figure 5. Posterior summaries of the coefficients of covariates for the probability of occupancy.

necessary tool for researchers and practitioners to analyse their single-species eDNA data and obtain reliable estimates of site-specific probabilities of species presence while accounting for false positive and false negative observation error.

Unlike previous R-packages for fitting multi-scale occupancy models that have been applied to eDNA data (Dorazio and Erickson 2018, Stratton et al. 2020), our implementation of the Griffin et al. (2020) model is novel in that it enables the estimation of false positive as well as false negative observation errors, both of which are known to be non-negligible in eDNA surveys. In addition, our RShiny app enables efficient Bayesian variable selection, which works well even when the number of predictors to be considered is large.

In terms of professional practice and conflicts involving the presence of protected species, a court of law may demand ‘proof’. Expressing presence/absence in terms of eDNA detection probabilities may therefore create uncertainties in making important decisions. As discussed by Griffiths et al. (2015), we urge practitioners to interpret ‘probability’ in terms of ‘risk’ level when it comes to decision-making. Ultimately, the level of certainty that is acceptable will therefore depend on the risk appetite for the decision concerned. Either way, ignoring uncertainty altogether is very high risk, as it is inherently

Table 3. Posterior summaries of the probabilities of a positive observation, true or false, in both stages of the survey.

Parameter	Mean (95% PCI)
θ_{11}	0.867 (0.802, 0.922)
θ_{10}	0.078 (0.004, 0.145)
p_{11}	0.862 (0.851, 0.872)
p_{10}	0.026 (0.023, 0.028)

Table 4. First row: Posterior probability of species absence conditional on $x=0, \dots, K$ positive qPCR replicates. Second row: Posterior probability of x positive qPCR replicated conditional on species presence. These correspond to baseline sites (all continuous covariates equal to 0 and all categorical covariates equal to their baseline level).

x	0	1	2	3	4	5	6	7	8	9	10	11	12
$1 - \psi(x)$	0.9766	0.9766	0.9766	0.9765	0.9462	0.4213	0.3683	0.3681	0.3681	0.3681	0.3681	0.3681	0.3681
$\psi(x)$	0.0976	0.0308	0.0045	0.0004	0.0001	0.0003	0.0023	0.0123	0.0477	0.1319	0.2465	0.2798	0.1458

present in all ecological sampling, whether it is estimated or not. We therefore recommend that practitioners embrace models that estimate detection probability and the risk of false positives and false negatives as a matter of course.

To cite RShiny app or acknowledge its use, cite this Software note as follows, substituting the version of the application that you used for ‘version 1.0’:

Diana, A. et al. 2021. An RShiny app for modelling environmental DNA data: accounting for false positive and false negative observation error. – *Ecography* 44: 1838–1844 (ver. 1.0).

Acknowledgements – We thank Natural England for collecting the data and making them available in an open access format. We also thank Dr Diana Cole for useful insights about the issue of model identifiability and Dr Pete Brotherton for advice and comments.

Funding – This work was funded by NERC project NE/T010045/1 ‘Integrating new statistical frameworks into eDNA survey and analysis at the landscape scale’.

Author contributions

Alex Diana: Conceptualization (equal); Formal analysis (equal); Software (lead); Writing – review and editing (supporting). **Eleni Matechou:** Conceptualization (equal); Formal analysis (equal); Project administration (lead); Supervision (equal); Writing – original draft (lead); Writing – review and editing (equal). **Jim E. Griffin:** Conceptualization (equal); Methodology (lead); Supervision (equal); Writing – review and editing (supporting). **Andrew S. Buxton:** Conceptualization (equal); Data curation (lead); Writing – review and editing (supporting). **Richard A. Griffiths:** Conceptualization (equal); Writing – review and editing (supporting).

Transparent Peer Review

The peer review history for this article is available at <<https://publons.com/publon/10.1111/ecog.05718>>.

Data availability statement

The example data set was collected by Natural England as part of a species distribution assessment project, and made available through the Natural England open data portal: <https://naturalengland-defra.opendata.arcgis.com/datasets/ffba3805a4d9439c95351ef7f26ab33c_0/data>.

References

- Attali, D. 2020. shinyjs: easily improve the user experience of your shiny apps in seconds. – R package ver. 2.0.0. <<https://cran.r-project.org/web/packages/shinyjs/index.html>>
- Bääth, R. 2018. beep: easily play notification sounds on any platform. – R package ver. 1.3. <<https://www.r-project.org/nosvn/pandoc/beep.html>>
- Barbieri, M. M. et al. 2004. Optimal predictive model selection. – *Ann. Stat.* 32: 870–897.

- Brown, P. J. et al. 1998. Multivariate Bayesian variable selection and prediction. – *J. R. Stat. Soc. B* 60: 627–641.
- Chang, W. 2018. shinythemes: themes for shiny. – R package ver. 1.1.2. <<https://cran.r-project.org/web/packages/shinythemes/index.html>>
- Chang, W. et al. 2020. shiny: web application framework for R. – R package ver. 1.5.0. <<https://cran.r-project.org/web/packages/shinyssloaders/index.html>>
- Chipman, H. et al. 2001. The practical implementation of Bayesian model selection. – *Lecture Notes-Monogr. Ser.* pp. 65–134.
- Cole, D. J. 2020. Parameter redundancy and identifiability. – CRC Press.
- Dorazio, R. M. and Erickson, R. A. 2018. ednaoccupancy: an R package for multiscale occupancy modelling of environmental DNA data. – *Mol. Ecol. Resour.* 18: 368–380.
- Ghosh, J. 2015. Bayesian model selection using the median probability model. – *Wiley Interdiscip. Rev. Comput. Stat.* 7: 185–193.
- Griffin, J. E. et al. 2020. Modelling environmental DNA data; Bayesian variable selection accounting for false positive and false negative errors. – *J. R. Stat. Soc. C* 69: 377–392.
- Griffiths, R. A. et al. 2015. Science, statistics and surveys: a herpetological perspective. – *J. Appl. Ecol.* 52: 1413.
- Guillera-Arroita, G et al. 2017. Dealing with false-positive and false-negative errors about species occurrence at multiple levels. – *Methods Ecol. Evol.* 8: 1081–1091.
- McClenaghan, B. et al. 2020. Validating metabarcoding-based biodiversity assessments with multi-species occupancy models: a case study using coastal marine eDNA. – *PLoS One* 15: e0224119.
- Oldham, R. et al. 2000. Evaluating the suitability of habitat for the great crested newt *Triturus cristatus*. – *Herpetol. J.* 10: 143–155.
- Plummer, M. et al. 2006. CODA: convergence diagnosis and output analysis for MCMC. – *R News* 6: 7–11.
- Polson, N. G. et al. 2013. Bayesian inference for logistic models using Pólya–Gamma latent variables. – *J. Am. Stat. Assoc.* 108: 1339–1349.
- Sali, A. and Attali, D. 2020. shinyssloaders: add loading animations to a ‘shiny’ output while it's recalculating. – R package ver. 1.0.0. <<https://cran.r-project.org/web/packages/shinyssloaders/index.html>>
- Stratton, C. et al. 2020. msocc: fit and analyse computationally efficient multi-scale occupancy models in R. – *Methods Ecol. Evol.* 11: 1113–1120.
- Thomsen, P. F. and Willerslev, E. 2015. Environmental DNA – an emerging tool in conservation for monitoring past and present biodiversity. – *Biol. Conserv.* 183: 4–18.
- Thomsen, P. F. et al. 2012. Monitoring endangered freshwater biodiversity using environmental DNA. – *Mol. Ecol.* 21: 2565–2573.
- Valentini, A. et al. 2016. Next-generation monitoring of aquatic biodiversity using environmental DNA metabar coding. – *Mol. Ecol.* 25: 929–942.
- Wickham, H. 2007. Reshaping data with the reshape package. – *J. Stat. Softw.* 21: 1–20.
- Wickham, H. 2016. ggplot2: elegant graphics for data analysis. – Springer-Verlag.
- Wickham, H. et al. 2020. dplyr: a grammar of data manipulation. – R package ver. 1.0.2. <<https://cran.r-project.org/web/packages/dplyr/index.html>>