
A four-step Bayesian workflow for improving ecological science

Keywords: big data; scientific workflow; data simulation; forecasting; null hypothesis testing

Example & Data: We provide an example of the workflow with complete code available eventually on its GitHub site, but for double-blind review we provide it here as a compressed set of files. The data used for the example is provided and full metadata on it is available via the Knowledge Network for Biocomplexity: [doi:10.5063/F12J69B2](https://doi.org/10.5063/F12J69B2).

Abstract

Growing anthropogenic pressures have increased the need for robust predictive models. Meeting this demand requires approaches that can handle bigger data to yield forecasts that capture the variability and underlying uncertainty of ecological systems. Bayesian models are especially adept at this and are growing in use in ecology. Yet many ecologists are not trained in current methods to build flexible robust models. Here we describe a broadly generalizable workflow for statistical analyses and show how it can enhance training in ecology. Building on the increasingly computational toolkit of many ecologists, this approach leverages simulation to integrate model building and testing for empirical data more fully with ecological theory. In turn this workflow can fit models that are more robust and well-suited to provide new ecological insights—allowing us to refine where to put resources for better estimates, better models, and better forecasts.

Introduction

In recent years, as ecologists have worked to develop global predictive models, they have developed ever larger datasets (Hampton et al., 2013). These bigger data, however, are also messier data. Such data generally requires a model of both the underlying biological processes and how the measurements were made. Some fields have long used these types of models (generally in fields focused on inferring population sizes for management, Muthukumarana et al., 2008; Zheng et al., 2007; Trijoulet et al., 2018; Strinella et al., 2020). Most, however, have not. This has left many researchers to try to adapt what they were trained in—traditional statistical methods (e.g. F and t tests) and a strong focus on null hypothesis testing—to increasingly complex datasets.

Yet many common statistical approaches do not align with ecology’s aims today. Beyond the reality that most traditional methods are fragile when used beyond the cleaner, simpler experiments these methods assume (e.g. spatial, temporal and phylogenetic correlations often violate independence assumptions), they will usually fail to produce robust, reproducible results. For example, an overly zealous focus on p -values has led to a replication crises in several fields, where results seem most likely the outcome of noisy data combined with a search for statistical significance through many models (effectively a garden of forking paths, Halsey et al., 2015; Loken and Gelman, 2017). Some model selection approaches, including new machine learning methods, try to avoid this by comparing across models, but may not generalize to provide useful forecasts. This is especially true when forecasts have to adapt to changes in the underlying biology.

Bayesian approaches provide a pathway to build powerful models that can transform how we understand our systems as largescale ecological data become increasingly available. Recognizing this, many in ecology are increasingly using Bayesian methods (Anderson et al., 2021). New algorithms (e.g. Hamiltonian Monte Carlo, Hoffman and Gelman, 2014; Betancourt, 2019) that have made fitting and implementing Bayesian models faster, more robust and—in many ways—easier (Carpenter et al., 2017). Fitting larger and sometimes more complex models,

however, presents challenges that are frequently not addressed in traditional ecological training. We suggest that many of these challenges can be overcome by approaching analyses through specific workflows (Betancourt, 2020; Grinsztajn et al., 2021; van de Schoot et al., 2021), which themselves are built on a process of how to do not just statistics, but how to do science (Box, 1976).

Such approaches move away from a focus on null hypothesis testing, towards estimating effect sizes, using models calibrated (see Table 1) and better understood through simulating data at multiple steps. Here we outline a simplified—but powerful—workflow that builds on new insights from statistics (Betancourt, 2020; van de Schoot et al., 2021) and the increasingly computational nature of ecology today. Our aim is to provide an approachable rubric for those new to fitting complex models, and is not intended to be a comprehensive overview (see ‘Next steps’ in the Supplement). We introduce our workflow assuming a Bayesian statistical framework (with an example in the supplement shown in R and Stan); however, it can be applied to other statistical inference methods.

A basic Bayesian workflow

Our workflow outlines what we consider the major steps for building bespoke models (Fig. 1). Such models can be fit by applying Bayes’ theorem, which generates a *posterior* distribution from a combination of a *likelihood* and a *prior* distribution (an initial uncertainty estimate derived from basic ecological knowledge), and using iterative algorithms (e.g. MCMC, Markov Chain Monte Carlo) that provide samples that can be used to extract information from the posterior distribution (for more, see *A brief review of statistical inference using Bayesian approaches* in the Supplement). Many of these steps will be familiar to statistical ecologists, but are often overlooked, whereas other steps may appear particular to Bayesian methods (e.g. prior predictive checks), but are actually useful for anyone—using Bayesian models or not—to challenge their models of how the world works. Parts of this workflow could be expanded as workflows in themselves, given other aims (see Supplement: Which workflow?).

Step 1: Develop your model(s)

We start the workflow with what can feel like the biggest step—build a model (or potentially, models) based on your aims. By developing a model designed for your biological question, data and aims, your statistical workflow naturally becomes a scientific workflow. You will more clearly see the assumptions and mechanisms in your model, which is especially valuable given how often our intuition of how models ‘work’ is wrong (Kokko, 2005). You likely already have a model, though it may be only verbal or conceptual. For this workflow, however, you’ll need to convert such models into mathematical versions (Servedio et al., 2014).

Though it can feel challenging at first, this step is best approached before you collect any data. A suite of resources for ‘generative’ or ‘narratively generative’ modeling can help (MacElreath, 2016; Betancourt, 2021*b*). As you start, ask lots of questions—and push yourself on your answers—about what you expect and what’s reasonable biologically from your model. As you do this, you’ll be generating your model—including its priors, which are important for Bayesian analysis. Assigning priors generally forces you to think about your model with regards to your study system, and interrogate what’s probable, possible or actually unreasonable—and can quickly disabuse users of prejudices regarding priors. For example, you may not think you have a prior on how sunlight affects plant growth, until you realize your ‘agnostic prior’ actually allows plants to grow hundreds of meters per day.

Step 2: Check your model on simulated data

Once you have your model and its priors jotted down, you need to write up your model in a particular modeling language and check it. As with all code: just because it runs, does not mean it does what you think it does. The worst errors often still permit code to run.

Test data (aka ‘simulated data’, or ‘fake data,’ etc.), and the skills required to build it, are central to this workflow. With ‘test data’ you simulate data from your model in such a way that you can use the resulting data to test if your model code is correct (i.e., you fix values for your model parameters, then test how well your model recovers them, see the Supplement

for an example). While there's no guarantee that inferences will always recover the parameter values you set, even when using the correct model, extreme disagreement is often an indicator that something is amiss in the implementation of the model. At the same time these simulation studies can help understand how often a model might lead to the correct inference (see Fig. 2). As you do this, you will also be calibrating your model—seeing how accurately and precisely it estimates parameters and under what conditions.

This very basic model checking step is uncommon for many ecologists, but critical in our view. If you can simulate data from your model, then you can powerfully—and easily—answer questions related to statistical power, what effect sizes are reasonable, and—most likely—have new insights into how your model suggests the world works, all before looking at any real data. Thus, this apparently simple programmatic task actually encapsulates a far deeper understanding of your model.

You can learn only so much, however, from data simulated from a particular parameter set. Simulation studies across multiple parameter sets allow you to investigate how robust your inferential performance might be. Prior predictive checks (Betancourt, 2021a; Wesner and Pomeranz, 2021; Winter and Depaoli, 2023) use the Bayesian prior model to set this scope of such simulations. For these, you draw values from your prior distribution and then explore how your model performs. Seeing how this influences your resulting output reveals the extent to which your model can capture known variation in your data, and gives insight into whether your model is capable of distinguishing among competing hypotheses.

Step 3: Run your model on your empirical data

The next step is to run the model—you've now evaluated, test-run and have ready to go—on your exciting new empirical data. Check diagnostics so you know it's running well and adjust until it is (this includes a suite of convergence and efficiency metrics that are well-discussed elsewhere, Betancourt, 2020; Gelman et al., 2020; van de Schoot et al., 2021; Gabry et al., 2019).

This is the step many ecologists skip straight to, ourselves included. It's easy to see the appeal: this is the inference step and where you might gain new ecological insights. Fitting new data to the model can feel like the moment when you'll learn something new. But, at least in our experience, this is not always the case. When we rush to this step, that first model we fit is often followed by another, and another—perhaps because one does not converge, or the results of another do not make immediate sense. And with the excitement of getting a model to run we can get distracted from what we are actually most interested in—the inference into ecology. Following this workflow can make this step much more satisfying. Here the benefits of the workflow may become excitedly apparent: you have estimates in useful units with uncertainty you can understand. You can use this information to draw new conclusions, design new experiments and more—but this is also a point to stop and check your model.

Step 4: Check your model on data simulated from your empirical model output (also known as posterior retrodictive checks)

Once you have your posterior based on your model and new empirical data, it's time to remember that it's wrong (as all models are) and ask how useful it is. This is where simulating from your model can be especially insightful. It will not only indicate that the model isn't adequately fitting the data but also can suggest what the problems might be. Using the parameter estimates from your posterior to simulate new data (Held et al., 2010; Gelman et al., 2000; Conn et al., 2018) lets you see how that new world compares to the observed data—called posterior retrodictive checks (or posterior predictive checks, Fig. 3).

Often here you may find big differences from your empirical data, and can start to generate hypotheses for why. For example, you may find patterns that suggest missing grouping factors (e.g. site or biome) through visual posterior retrodictive checks, or you may quickly realize your model predicts impossible numbers for your biological reality. You may begin to see inadequacies in your model, or even potentially your data. This is one of the main benefits of the workflow: models don't fail silently, they fail with a wealth of context that helps to generate new models

and experiments.

Feedbacks & workflows

A key feature of this workflow is that it can be iterated. If you find that you want to tweak your model then you return to the beginning, adjust your model, and repeat the rest of the workflow. In this way, fitting multiple models is encouraged, but this is distinct from the quest for a minimum adequate model or one ‘best’ fit. Feedbacks in this workflow are focused far more on what is biologically reasonable, and understanding the utility—and limits—of inference from your data for your model. And there are big benefits to it.

How this workflow changed our science

Before this workflow, not all of us commonly discussed the values that parameters in our model took—things like the slope and intercept (two common model parameters) were sometimes reported, but we did not know them as well as we knew whether the p -value for the slope was < 0.05 . This changes quickly when you need to build simulated data (Step 2). For example, when modeling phenological events (observations of biological events on numbered days within the calendar year: 1-365 most years) it is not uncommon to find seemingly-reasonable models generating predictions of events on the non-existent calendar day of 1000.

A closer inspection of our parameters also taught us a lot about identifiability and nonidentifiability, which refers to when all parameters in a model can—or cannot—be uniquely identified with infinite data, and a statistical kin: degeneracy (see Table 1). Degeneracy concerns the kinds of complex uncertainties that can arise from finite data sets (Gelman and Hill, 2009), and something we have often found in Steps 2-3 of our workflow. Nonidentifiability and degeneracy can come up in many ways in ecology, and make us think we understand processes we do not. We never noticed them before using this workflow, but since then we have realized (especially in steps 1-2) lots of places for nonidentifiability and degeneracies to live—and we have adjusted

how we collect data and interpret results because of it. For example, we have found fitting both site and species in a model with highly imbalanced data or trying to estimate interaction terms with low sample sizes (Gelman and Hill, 2020, for more details) leads to degenerate models, but there’s often no warning in packages to tell us this.

How this workflow intersects with ecological training

This four-step workflow is a simplified version of the current best practices for Bayesian model fitting (Betancourt, 2020; van de Schoot et al., 2021), but many of the skills required are not part of traditional ecological training. Writing out the math behind most statistical models enough to complete Steps 1-2 bleeds into the skillset usually reserved for those working on theory, where coding and simulating from a model are common tasks. In contrast field, lab and otherwise empirical-data based ecologists often fit models they could not simulate data from. This dichotomy seems short-sighted in our current era of bigger, messier data and greater computational methods poised to handle such messy data. Further, the increasingly computational toolkit of the modern ecologist makes it easier to bridge the gap between ecological models and their underlying math.

We argue training in simulating data as part of an organized workflow could speed progress in ecology and is possible given the current skillset of many ecologists. A reasonably competent coder could easily simulate data under a complex model that they might not have the mathematical expertise to solve analytically—if doing so was part of their training and the workflows they regularly use.

Advances in developing Bayesian workflows have come alongside improved algorithms, visualizations (e.g. Betancourt, 2020; van de Schoot et al., 2021; Gabry et al., 2019), perspectives on priors (Gelman et al., 2014; Gelman and Hill, 2020; Betancourt, 2021*a*) and hierarchical approaches that could also improve training. For example, new work shows that prior predictive checks provide a more powerful and intuitive way to understand how priors work within a particular model (Betancourt, 2021*a*), compared to past approaches. Similarly, traditional

ecological training in hierarchical models still often refers to grouping factors (such as species or individual) as ‘random effects,’ which is misleading, imprecise and thus no longer recommended (Gelman and Hill, 2009).

These new best practices have gained traction at the same time many fields have recognized that p -values, and null hypothesis testing in general, are easily misleading. Small sample sizes, lack of routine reporting of interpretable effect sizes, fitting of many models without adequate explanation, poor data and code reporting habits all increase the chance of finding ‘significance’ at a level of ≤ 0.05 (Halsey et al., 2015; Loken and Gelman, 2017). Small sample sizes alongside a tendency to fit complicated models with multiple interactions makes ecological research vulnerable to these problems. The answer to this, however, is not to make p -values smaller (Halsey et al., 2015; Colquhoun, 2017), nor is it Bayesian approaches. The answer is training in workflows designed for careful model building, model fitting and model interrogation informed by ecological theory and understanding—including the one we outline here.

References

- Anderson, S. C., P. R. Elsen, B. B. Hughes, R. K. Tonietto, M. C. Bletz, D. A. Gill, M. A. Holgersson, S. E. Kuebbing, C. McDonough MacKenzie, M. H. Meek, et al. 2021. Trends in ecology and conservation over eight decades. *Frontiers in Ecology and the Environment* 19:274–282.
- Betancourt, M. 2019. The Convergence of Markov Chain Monte Carlo Methods: From the Metropolis Method to Hamiltonian Monte Carlo. *Annalen der physik* 531.
- . 2020. Towards A Principled Bayesian Workflow. https://betanalpha.github.io/assets/case_studies/principled_bayesian_workflow.html.
- . 2021a. Prior modeling. https://betanalpha.github.io/assets/case_studies/prior_modeling.html.
- . 2021b. (what’s the probabilistic story) modeling glory? https://betanalpha.github.io/assets/case_studies/generative_modeling.html.
- Box, G. E. 1976. Science and statistics. *Journal of the American Statistical Association* pages 791–799.
- Carpenter, B., A. Gelman, M. Hoffman, D. Lee, B. Goodrich, M. Betancourt, M. A. Brubaker, J. Guo, P. Li, and R. Allen. 2017. Stan: A probabilistic programming language. *Journal of Statistical Software* 76:10.18637/jss.v076.i01.
- Colquhoun, D. 2017. The reproducibility of research and the misinterpretation of p-values. *Royal Society Open Science* 4.
- Conn, P. B., D. S. Johnson, P. J. Williams, S. R. Melin, and M. B. Hooten. 2018. A guide to Bayesian model checking for ecologists. *Ecological Monographs* 88:526–542.
- Gabry, J., D. Simpson, A. Vehtari, M. Betancourt, and A. Gelman. 2019. Visualization in

bayesian workflow. *Journal of the Royal Statistical Society Series a-Statistics in Society*
182:389–402.

Gelman, A., J. B. Carlin, H. S. Stern, D. B. Dunson, A. Vehtari, and D. B. Rubin. 2014.
Bayesian Data Analysis. 3rd ed. CRC Press, New York.

Gelman, A., Y. Goegebeur, F. Tuerlinckx, and I. Van Mechelen. 2000. Diagnostic checks for
discrete data regression models using posterior predictive simulations. *Journal of the Royal*
Statistical Society Series C-Applied Statistics 49:247–268.

Gelman, A., and J. Hill. 2009. *Data Analysis Using Regression and Multilevel/Hierarchical*
Models. Cambridge, New York.

———. 2020. *Regression and Other Stories*. Cambridge University Press.

Gelman, A., A. Vehtari, D. Simpson, C. C. Margossian, B. Carpenter, Y. Yao, L. Kennedy,
J. Gabry, P.-C. Bürkner, and M. Modrák. 2020. Bayesian workflow. *arXiv*.

Grinsztajn, L., E. Semenova, C. C. Margossian, and J. Riou. 2021. Bayesian workflow for disease
transmission modeling in stan. *Statistics in Medicine* 40:6209–6234.

Halsey, L. G., D. Curran-Everett, S. L. Vowler, and G. B. Drummond. 2015. The fickle p value
generates irreproducible results. *Nature Methods* 12:179–185.

Hampton, S. E., C. A. Strasser, J. J. Tewksbury, W. K. Gram, A. E. Budden, A. L. Batcheller,
C. S. Duke, and J. H. Porter. 2013. Big data and the future of ecology. *Frontiers in Ecology*
and the Environment 11:156–162.

Held, L., B. Schroedle, and H. Rue. 2010. Posterior and Cross-validatory Predictive Checks: A
Comparison of MCMC and INLA. Pages 91–110 *in* T. Kneib and G. Tutz, eds. *Statistical*
Modelling and Regression Structures: Festschrift in Honour of Ludwig Fahrmeir.

Hoffman, M. D., and A. Gelman. 2014. The No-U-Turn Sampler: Adaptively Setting Path
Lengths in Hamiltonian Monte Carlo. *Journal of Machine Learning Research* 15:1593–1623.

Kokko, H. 2005. Useful ways of being wrong. *Journal of evolutionary biology* 18:1155–1157.

-
- Loken, E., and A. Gelman. 2017. Measurement error and the replication crisis. *Science* 355:584–585.
- MacElreath, R. 2016. *Statistical Rethinking*, vol. 469 pp. CRC Press, New York.
- Muthukumarana, S., C. J. Schwarz, and T. B. Swartz. 2008. Bayesian analysis of mark-recapture data with travel time-dependent survival probabilities. *Canadian Journal of Statistics* 36:5–21.
- Servedio, M. R., Y. Brandvain, S. Dhole, C. L. Fitzpatrick, E. E. Goldberg, C. A. Stern, J. Van Cleve, and D. J. Yeh. 2014. Not just a theory—the utility of mathematical models in evolutionary biology. *PLoS biology* 12:e1002017.
- Strinella, E., D. Scridel, M. Brambilla, C. Schano, and F. Korner-Nievergelt. 2020. Potential sex-dependent effects of weather on apparent survival of a high-elevation specialist. *Scientific Reports* 10:8386.
- Trijoulet, V., S. J. Holmes, and R. M. Cook. 2018. Grey seal predation mortality on three depleted stocks in the West of Scotland: What are the implications for stock assessments? *Canadian Journal of Fisheries and Aquatic Sciences* 75:723–732.
- van de Schoot, R., S. Depaoli, R. King, B. Kramer, K. Maertens, M. C. Tadesse, M. Vannucci, A. Gelman, D. Veen, J. Willemsen, and C. Yau. 2021. Bayesian statistics and modelling. *Nature Reviews Methods Primers* 1.
- Wesner, J. S., and J. P. F. Pomeranz. 2021. Choosing priors in Bayesian ecological models by simulating from the prior predictive distribution. *Ecosphere* 12.
- Winter, S. D. D., and S. Depaoli. 2023. Illustrating the value of prior predictive checking for bayesian structural equation modeling. *Structural Equation Modeling-a multidisciplinary journal* .
- Zheng, C., O. Ovaskainen, M. Saastamoinen, and I. Hanski. 2007. Age-dependent survival analyzed with Bayesian models of mark-recapture data. *Ecology* 88:1970–1976.

Table 1: Glossary: We provide below simplified definitions of the major terms we use (many of these terms, such as calibration, may be used differently depending on the particular literature).

<i>Term</i>	<i>Definition</i>
calibration	analyzing how often an estimate is close to the true value over an ensemble of hypothetical observations. An exact calibration would requires simulating from the true data generating process which is impossible in practice. We can, however, calibrate to data simulated from the configurations of models we plan use to fit to our data (<i>Steps 1-2</i>) so we understand the models better, including their limits given data similar to ours. We emphasize simulations to calibrate model behaviors consistent with our ecological systems and understanding (e.g. working within a limited set of parameter ranges through prior predictive checks). In contrast to this approach, frequentist method are calibrated against all possible behaviors, which is not only impractical for complicated models it's also irrelevant given that the most extreme behaviors are unlikely to manifest in reality.
degeneracy	complex uncertainties that come from a mix of sources, including, non-identified models and cases where the data cannot well inform model parameters. When the data are not informing the parameters that we care about, this highlights a measurement issue. Identifying these problems in simulation studies can highlight when we need a better experimental design (e.g. sampling for more overlapping species across sites, or changing what we measure, etc.).
non-identifiability	when all parameters in a model cannot be uniquely identified with infinite data
prior	an distribution of reasonable values for a parameter based on fundamental biological and ecological understanding, previous research, or other sources
statistical model	Mathematical approximations of the true data generating process labeled with numerical parameters. Evaluating a statistical model on observed data gives

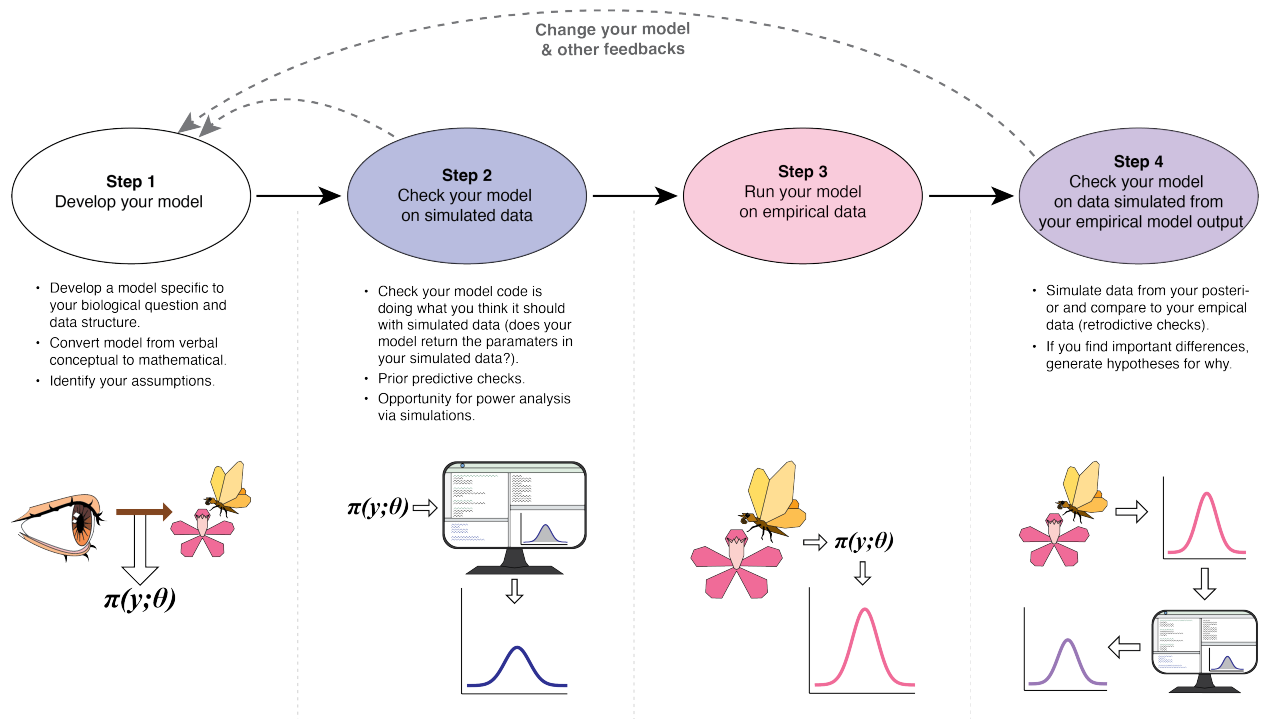


Figure 1: The four-step iterative workflow we outline can help design models for specific ecological questions, data and aims—which makes this a statistical workflow that can naturally become a scientific workflow. It makes the step that many of us focus on—running your model on your empirical data (Step 3)—far more straightforward and insightful by using simulations both before (Step 2) and after (Step 4) it to better understand the model and data together.

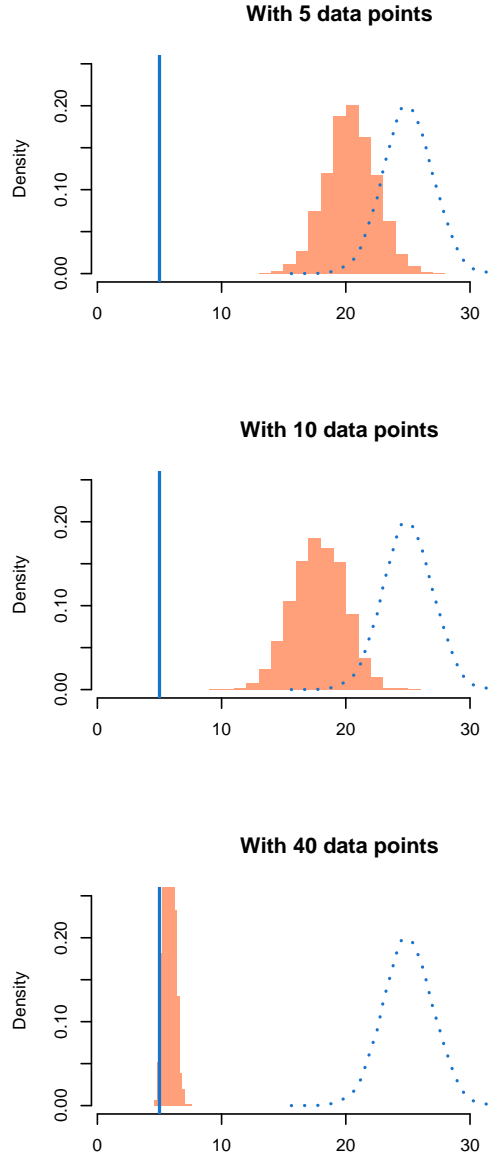


Figure 2: A simple example of how to use simulated data to understand calibration issues in a mis-specified model. Here we know the true model underlying the data is $y = \alpha + \text{normal}(0, \sigma)$ where α is 5 (shown as blue vertical line) and σ is 2. The model, however, is mis-specified by a prior for α of $\text{normal}(25, 2)$ (dashed blue line), resulting in a posterior (salmon-colored histogram) not centered on the true value. In our experience it is quite rare to have a prior informed by ecological knowledge be so far off, but this is an example. How mis-calibrated the model will be depends on the data: we show examples with a sample size (N) of 5, 10 and 40 data points. In practice these studies would allow us to determine how much data we would need to be robust to suspect prior models.

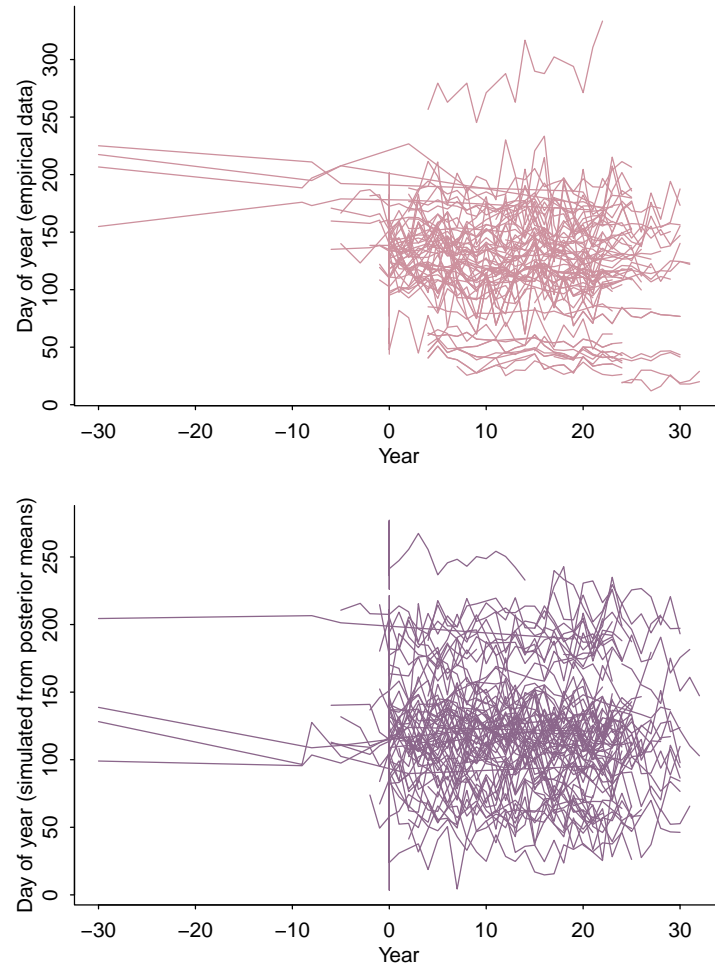


Figure 3: Example of a single retrodictive check from time-series data of phenological events over time. The raw data (top, pink) looks similar to one simulated dataset (bottom, purple), based on existing species number, their respective x data, and simulating from the parameters for each species. See ‘An example workflow’ in the Supplement for more details.