

## A guide to Bayesian model selection for ecologists

M. B. HOOTEN<sup>1,2,3,4,7</sup> AND N. T. HOBBS<sup>4,5,6</sup>

<sup>1</sup>*U.S. Geological Survey, Colorado Cooperative Fish and Wildlife Research Unit, Colorado State University, Fort Collins, Colorado 80523-1484 USA*

<sup>2</sup>*Department of Fish, Wildlife, and Conservation Biology, Colorado State University, Fort Collins, Colorado 80523-1484 USA*

<sup>3</sup>*Department of Statistics, Colorado State University, Fort Collins, Colorado 80523-1484 USA*

<sup>4</sup>*Graduate Degree Program in Ecology, Colorado State University, Fort Collins, Colorado 80523-1484 USA*

<sup>5</sup>*Department of Ecosystem Science and Sustainability, Colorado State University, Fort Collins, Colorado 80523-1484 USA*

<sup>6</sup>*Natural Resource Ecology Laboratory, Colorado State University, Fort Collins, Colorado 80523-1484 USA*

*Abstract.* The steady upward trend in the use of model selection and Bayesian methods in ecological research has made it clear that both approaches to inference are important for modern analysis of models and data. However, in teaching Bayesian methods and in working with our research colleagues, we have noticed a general dissatisfaction with the available literature on Bayesian model selection and multimodel inference. Students and researchers new to Bayesian methods quickly find that the published advice on model selection is often preferential in its treatment of options for analysis, frequently advocating one particular method above others. The recent appearance of many articles and textbooks on Bayesian modeling has provided welcome background on relevant approaches to model selection in the Bayesian framework, but most of these are either very narrowly focused in scope or inaccessible to ecologists. Moreover, the methodological details of Bayesian model selection approaches are spread thinly throughout the literature, appearing in journals from many different fields. Our aim with this guide is to condense the large body of literature on Bayesian approaches to model selection and multimodel inference and present it specifically for quantitative ecologists as neutrally as possible. We also bring to light a few important and fundamental concepts relating directly to model selection that seem to have gone unnoticed in the ecological literature. Throughout, we provide only a minimal discussion of philosophy, preferring instead to examine the breadth of approaches as well as their practical advantages and disadvantages. This guide serves as a reference for ecologists using Bayesian methods, so that they can better understand their options and can make an informed choice that is best aligned with their goals for inference.

*Key words:* Akaike information criterion; Bayes factors; cross-validation; deviance information criterion; model averaging; multi-model inference; regularization; shrinkage.

### INTRODUCTION

Model selection and Bayesian statistics have become increasingly important tools in the field of ecology (Johnson and Omland 2004, Clark 2005, Cressie et al. 2009, Hobbs 2009). Despite an upward trend in the use of model selection and Bayesian methods in ecological research, the intersection of these two frameworks for inference has been minimal in the literature (Fig. 1). The guidance provided about model selection in the Bayesian statistical literature is unbalanced and lacks cohesion.

The theory and protocol for implementing a variety of Bayesian model selection methods seem much less tangible than the information criterion approaches for maximum likelihood we have grown accustomed to in ecology. Thus, we are at a critical juncture in our field. Do we use newer statistical technology while potentially foregoing model selection because it is too complicated, or do we use more familiar statistical methods at the potential risk of letting our choice of selection procedure dictate what scientific questions we can answer with our model(s)? An awareness of available model comparison approaches in the Bayesian framework can help the ecologist choose and apply the method that is most suited to their goals for inference.

Manuscript received 7 April 2014; accepted 12 May 2014.  
Corresponding Editor: A. M. Ellison.

<sup>7</sup> E-mail: mevin.hooten@colostate.edu

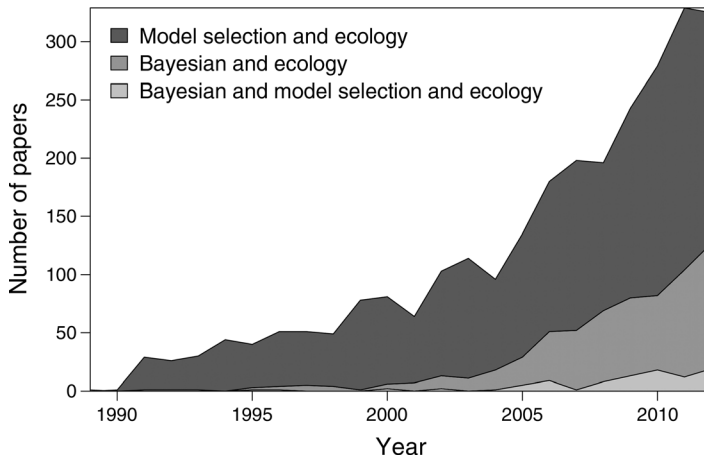


FIG. 1. The results of a Web of Science search in number of articles per search string for each of the past 25 years.

#### *Preliminary assumptions and notation*

Our primary focus is on providing a comprehensive description of available methods for Bayesian model selection and multimodel inference that is accessible to ecologists. For a discussion of the philosophical arguments pertaining to model selection and multimodel inference we refer the interested reader to several excellent sources, including Gelman and Shalizi (2012) and J. M. Ver Hoef and P. L. Boveng (*unpublished manuscript*), who discuss when and why one should use model selection methods. In this exposition, we assume the reader is familiar with the philosophical underpinnings and has already decided that they (1) seek Bayesian statistical inference, (2) would like to compare models for the purpose of improving that inference, and (3) have already verified the model assumptions for their particular data set. This last item is critical because if the model assumptions are not met, the resulting statistical inference (including predictions and prediction uncertainty) rests on a house of cards. Reliable inference requires checking the assumptions of our models. For further details on model checking, including the evaluation of goodness of fit and posterior predictive  $P$  values, see Gelman et al. (2014a).

We also assume the reader has broad familiarity with statistical methods including least squares and maximum likelihood, as well as a basic understanding of Bayesian model building and algorithms for implementation (e.g., Markov chain Monte Carlo). Gotelli and Ellison (2012) and Bolker (2008) provide excellent background on contemporary ecological statistics, and from a Bayesian perspective see Clark (2007), Royle and Dorazio (2008), Link and Barker (2010), and Hobbs and Hooten (*in press*).

We make frequent use of matrix notation and linear algebra (to avoid excessive summation notation) throughout this guide, but readers unfamiliar with these concepts will be able to glean the big-picture concepts and connections from our descriptions. In particular, we use a common Bayesian square bracket notation  $[a|b]$

(courtesy of Gelfand and Smith 1990) to represent probability distributions, in this case, the distribution of variable  $a$  given variable  $b$ . We also make occasional use of the probability notation  $P(c)$  to denote the probability of item  $c$ . For matrix notation, we use a standard form where matrices and vectors are bold, with matrices uppercase (e.g.,  $\mathbf{X}$ ) and vectors lowercase (e.g.,  $\mathbf{x}$ ). Transpose matrix and vector is denoted by the prime symbol (e.g.,  $\mathbf{x}'$ ). We use  $\theta$  generically to denote a set of model parameters, and  $\mathbf{y}$  to denote a data set, typically composed of response variables. Finally, we have defined several commonly used terms in the model selection and Bayesian literature in Table 1 to aid those readers less familiar with the subject.

#### *Overview of topics*

In this guide, we present a wealth of available perspectives on Bayesian multimodel inference and model selection. It may come as a surprise that there are many options for model selection and multimodel inference, each with its own strengths and weaknesses. It is our view that ecologists need the ability to distinguish among methods more than they need a strict set of rules to follow in how to proceed with model selection. We use the term “guide” here (in the same sense as a field guide for birds) because we have made an effort to be thorough and to remain unaffiliated in our description of these methods. Our guide is intended to be used as a conceptual aid; ecologists can use it to learn about the variety of options available and can decide how each fits in with their own research goals. For illustration, we implement several specific methods (all computer code is available in the Supplement). However, as space does not allow us to provide specific examples of computational algorithms for every approach, we have made an effort to provide the reader with numerous references they can consult to implement these methods in the statistical software of their choice.

This paper is organized as follows. We begin by highlighting a few of most important and sometimes lesser known take-home messages concerning model

TABLE 1. Glossary of terms and definitions.

Term	Definition
AIC	Akaike's information criterion, a within-sample non-Bayesian score for prediction
Bayes factor	the ratio of marginal data distributions pertaining to two models
BIC	Bayesian (Schwartz) information criterion, a within-sample non-Bayesian score for model averaging
CPO	conditional predictive ordinate, a within-sample score for leverage
Cross-validation	the iterative use of within-sample data to validate models in terms of out-of-sample predictive ability
DIC	deviance information criterion, a within-sample quasi-Bayesian score for prediction
Effective number of parameters	$p_D$ , a measure of model complexity as a penalty in Bayesian information criteria
Empirical Bayesian	the use of within-sample data to inform Bayesian model components such as priors
Out-of-sample data	an auxiliary set of data that are used for model comparison
Posterior predictive loss	an approach for scoring models based on decision theory
Regularization	constraining a statistical optimization problem (i.e., penalization or shrinkage)
Regulator	constraint, optimism, penalty, or prior
Score	a function used to evaluate models numerically, usually in terms of predictive ability
WAIC	Watanabe-Akaike information criterion, a within-sample fully-Bayesian score for prediction
Within-sample data	response data typically used to fit a model, but also to calculate information criteria

selection. This prelude serves as an overview containing big picture connections between the methods we describe subsequently. We then introduce a specific Bayesian ecological model as a case example. We refer to this example throughout to illustrate differences among alternative approaches. We describe Bayesian model averaging, for use when the goal of the researcher is to make inferences from more than one model. We treat out-of-sample validation, the gold standard for model selection based on predictive ability. We then turn to a topic that applies broadly across Bayesian and non-Bayesian statistics, the process of regularization, which we feel is essential to understanding the subsequent material on information criteria. We then cover model-based methods for model selection. In the penultimate section, we provide specific guidance on matching alternative methods to inferential goals. As a visual aid to the flow of the manuscript, we show section topics and sub-topics in Fig. 2, providing an overview for the relationships among ideas and methods that we describe throughout the paper.

### Highlights

While preparing this guide, we experienced several epiphanies ourselves that had not occurred to us previously. We discovered that most of these findings have existed in the literature for quite some time (a decade, at least), but had not been brought together in a way that supports a solid understanding and intuition about model selection. Among the most important of our own epiphanies were the following (see Table 1 for definitions of terms):

- 1) There is no general consensus among statisticians on the topic of model selection.
- 2) Multimodel inference can be thought of from many different perspectives, including model averaging. Thus, we use the phrase “model selection” somewhat generically (including model comparison and multimodel inference) because many of the methods we describe inherently consider multiple models (sometimes infinitely many), but aren't considered to be model averaging in the conventional sense.
- 3) Much of the statistical community relies heavily on out-of-sample model comparison approaches, yet, in ecology, we primarily favor information criterion approaches that avoid the use of out-of-sample data for model evaluation. Despite the potential advantages for model selection, out-of-sample methods have been largely ignored by ecologists because they may require additional data beyond what was already collected in the study and historically were very computationally intensive to implement.
- 4) Cross-validation is a hybrid approach containing both out-of-sample and within-sample aspects. From a Bayesian perspective, cross-validation for model selection is considered to be an empirical Bayesian method and can be incredibly helpful for model selection.
- 5) Neither AIC nor BIC are appropriate for Bayesian model averaging in all situations. Both AIC and BIC were designed to be used with maximum likelihood estimates and make fairly strong assumptions about a priori model probabilities. Whereas AIC excels at finding good predictive models, BIC was developed mainly for model averaging purposes and is good for small sets of well-justified models.
- 6) DIC and AIC often yield quite similar results for model selection with certain classes of models, however, DIC is not ideal for all classes of models (e.g., mixture models). No theoretical justification exists in the literature for the use of DIC in model averaging. Furthermore, DIC is not a fully Bayesian model comparison criterion.
- 7) A truly Bayesian information criterion seems to have just been discovered (i.e., WAIC), but in actuality went unnoticed for more than a decade. WAIC resolves many of the issues with DIC, but

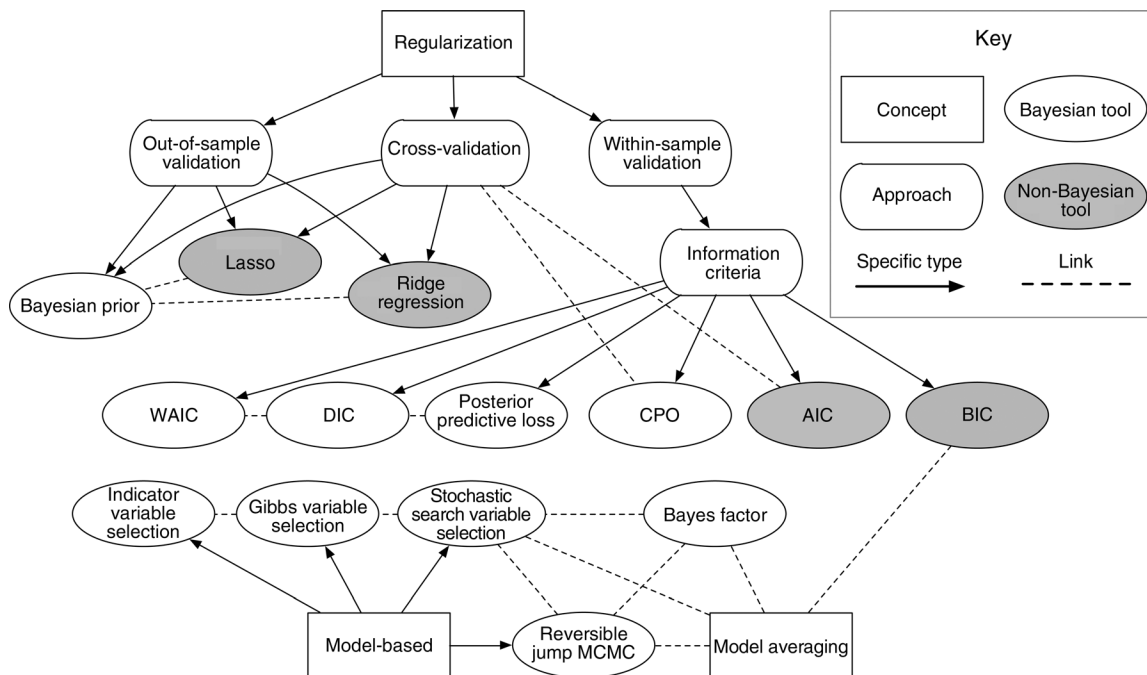


FIG. 2. Overview of topics treated in this guide. These topics are grouped by their linkages to the main model selection and multimodel inference themes. Boxes represent overarching concepts, rounded boxes represent certain approaches that fall under those concepts, and ovals correspond to specific tools (gray indicates tools that are not clearly Bayesian). Arrows indicate specific types of approaches and tools that fall under the broader concepts, whereas dashed lines represent links among items if certain assumptions hold (e.g., Bayesian information criterion [BIC] can be used for model averaging if parameters can easily be counted, priors are vague, and posterior modes are used as point estimates for parameters). Abbreviations are explained in Table 1.

also seems to have a critical weakness for some models.

- 8) Regularization is an umbrella concept that spans nearly all topics in model selection. When statistical optimization problems are written as regularization expressions, it becomes clear that AIC, BIC, DIC, WAIC, posterior predictive loss, ridge regression, and Lasso all fall under the same umbrella. Moreover, regularization itself has an inherently Bayesian justification. It explicitly constrains model parameters in the same way a Bayesian prior does. Thus, model selection is similar to using a strong prior, at least in spirit.
- 9) The Bayesian framework allows one to actually build parametric mechanisms into models that perform model selection (e.g., stochastic search variable selection and reversible jump Markov chain Monte Carlo [MCMC]). We refer to these as model-based model selection approaches. They can be viewed as a combination of model selection and multimodel inference.

*An exemplar: the hierarchical Bayesian occupancy model*

Mixture models, especially zero-inflated models, comprise an important class of statistical tools in contemporary ecological research. In particular, occupancy and capture–recapture models are very commonly

used in the field of wildlife ecology (Royle and Dorazio 2008). We consider the hierarchical occupancy model as a prototypical Bayesian ecological model. The Bayesian occupancy model presents challenges for traditional model comparison methods, thus, we introduce the model here and refer back to it later to demonstrate several approaches for model selection and multimodel inference.

In essence, the occupancy model is simply a binary regression model with binary measurement error. In its application, the occupancy model can be used to learn about the true presence or absence of a species and the niche-related features of the sites while accounting for imperfect detection (MacKenzie et al. 2006). The basic occupancy model, presented for ecologists, was described by MacKenzie et al. (2002) and included implementation details from a maximum likelihood perspective. More recently, occupancy models have been extended to model temporal dynamics (e.g., MacKenzie et al. 2003), spatial autocorrelation (e.g., Johnson et al. 2013), and community dependence (e.g., Dorazio et al. 2010).

Hierarchically, a simple occupancy model with homogeneous detection probability and heterogeneous occupancy probabilities can be written as a zero-inflated binomial data model (with detection probability  $p$ ) that depends on a latent Bernoulli process ( $z_i$ , presence or absence) that varies among sites ( $i = 1, \dots, n$ ) according



to probability  $\psi_i$ . The response data,  $y_i$ , are a sum of the binary detection history for each site over a set of visits or occasions ( $J_i$ ); that is,  $y_i = \sum_{j=1}^{J_i} y_{ij}$ , where  $y_{ij}$  are binary detection observations for site  $i$  on survey occasion  $j$ . On each occasion, the species is detected (i.e.,  $y_{ij} = 1$ ) with probability  $p$  if it is truly present, otherwise it is recorded as not detected (i.e.,  $y_{ij} = 0$ ). For simplicity, we have used a specification of the occupancy model that assumes a homogeneous detection probability  $p$  and conditional independence for detection on each site visit  $j = 1, \dots, J_i$ . These assumptions can be relaxed by allowing for variation in detection as well as occupancy probability.

The logit link,  $\log(\psi_i/(1 - \psi_i))$ , is most commonly used function relating occupancy probability  $\psi_i$  to a set of site-level covariates  $\mathbf{x}_i$ , however there can be computational advantages to using other link functions such as the probit (Hooten et al. 2003, Dorazio and Rodriguez 2012, Johnson et al. 2013). The probit link function allows us to reparameterize the model using a set of auxiliary variables  $v_i$  that describe a continuous latent process representing occupancy probability (Albert and Chib 1990). The probit occupancy model is specified hierarchically as

$$y_i \sim \begin{cases} 0 & \text{if } z_i = 0 \\ \text{Binom}(J_i, p) & \text{if } z_i = 1 \end{cases} \quad (1)$$

$$z_i \sim \begin{cases} 0 & \text{if } v_i \leq 0 \\ 1 & \text{if } v_i > 0 \end{cases} \quad (2)$$

$$v_i \sim N(\beta_0 + \mathbf{x}_i' \boldsymbol{\beta}, 1) \quad (3)$$

$$p \sim \text{Beta}(1, 1), \quad (4)$$

$$\beta_0 \sim N(\mu_0, \sigma_0^2), \quad (5)$$

$$\boldsymbol{\beta} \sim N(\boldsymbol{\mu}, \sigma_{\boldsymbol{\beta}}^2 \mathbf{I}), \quad (6)$$

where the probit link function itself (i.e.,  $\Phi$ , the standard normal cumulative distribution function) only comes into play when we condition  $z_i$  on the regression coefficients  $\beta_0$  and  $\boldsymbol{\beta}$  (e.g.,  $\boldsymbol{\beta} = (\beta_1, \beta_2)'$ ) directly; then we obtain  $z_i \sim \text{Bernoulli}(\Phi(\beta_0 + \mathbf{x}_i' \boldsymbol{\beta}))$ . The advantages of this probit occupancy model are primarily computational. The implicit probit link function allows us to create a fully Gibbs MCMC algorithm that requires no Metropolis-Hastings updates or tuning (Dorazio and Rodriguez 2012, Johnson et al. 2013). We use the probit occupancy model presented in Eqs. 1–6 as a basis for demonstrating the model selection procedures that follow, making modifications to it as needed.

#### MODEL AVERAGING

From here forward, assume that we are dealing with a set of models  $\mathcal{M} = \{M_1, \dots, M_l, \dots, M_L\}$  (where  $L$  is the total number of models) that are built using expert scientific judgement and are not obviously inappropriate

in terms of assumptions. Model averaging allows us to combine the strengths of several models for improved inference. It has been argued (e.g., Kass and Raftery 1995, Link and Barker 2006) that Bayesian model averaging (BMA) is the proper way to obtain multi-model inference under the Bayesian statistical paradigm because it provides a valid probability-based mechanism for considering multiple models in the presence of process and parameter uncertainty. Hoeting et al. (1999) provided an excellent overview of BMA, complete with implementation details for selected model classes.

An important and often overlooked aspect of model averaging is that BMA was not designed as a method for model selection, but rather as a method for combining posterior distributions. Whereas many of the methods in the following sections are based heavily on finding models that excel at out-of-sample predictive performance (e.g., AIC and DIC), BMA is intended for within-sample model combination. Thus, in what follows, we provide some insight about how BMA fits into the larger suite of model selection methods and refer the interested reader to the literature cited herein for details.

At the heart of BMA is the average posterior distribution of a quantity of interest ( $g \equiv g(\boldsymbol{\theta}, \tilde{\mathbf{y}})$ ), typically a function of either an unknown parameter or set of data or both)

$$[g | \mathbf{y}] = \sum_{l=1}^L [g | \mathbf{y}, M_l] P(M_l | \mathbf{y}) \quad (7)$$

where  $[g | \mathbf{y}, M_l]$  is the posterior distribution of  $g$  under individual model  $M_l$  and  $P(M_l | \mathbf{y})$  is the posterior probability of model  $M_l$ . The posterior model probability  $P(M_l | \mathbf{y})$  is the workhorse of the BMA procedure, providing the weight of evidence in the average Eq. 7 for one model over others. Thus, we have a natural and proper Bayesian framework for multimodel inference as long as we can find the required quantities in Eq. 7. Furthermore, BMA performed on a set of models  $\mathcal{M}$  yields better inference about  $g$  than any one of the models alone (Madigan and Raftery 1994), thus we have a compelling reason to use it.

#### *The utility of the marginal data distribution*

Recall the classical expression for Bayes rule assuming a single model

$$[\boldsymbol{\theta} | \mathbf{y}] = \frac{[\mathbf{y} | \boldsymbol{\theta}][\boldsymbol{\theta}]}{[\mathbf{y}]} \quad (8)$$

where  $[\boldsymbol{\theta}]$  is the prior distribution for the parameters. The denominator  $[\mathbf{y}]$ , which we typically avoid finding analytically, corresponds to the aforementioned marginal data distribution for the given model; it will be large for the same set of data if the model represents them well and small if it doesn't. The marginal data distribution  $[\mathbf{y}]$  is a natural model discrimination measure by itself and is

fundamental in computing the posterior model probabilities  $P(M_l | \mathbf{y})$ . To show this, we generalize the notation to include information concerning the individual model each  $[\mathbf{y}]$  is associated with. Therefore, let  $[\mathbf{y} | M_l]$  be the marginal data distribution for model  $l$ . Then, the posterior model probability can be written as

$$P(M_l | \mathbf{y}) = \frac{[\mathbf{y} | M_l]P(M_l)}{\sum_{j=1}^L [\mathbf{y} | M_j]P(M_j)} \quad (9)$$

where  $P(M_l)$  is the assumed prior model probability, which is commonly set to  $1/L$ . The use of equal prior model probabilities explicitly assumes that there may be no reason to prefer one model over another. The alternative is to set the  $P(M_l)$  such that they represent an a priori understanding of differences among model importance as long as the sum of prior model probabilities over all models in the set equals 1. To obtain the necessary marginal data distribution for model  $l$  we need to integrate over the parameters in the joint distribution of the data  $\mathbf{y}$ , the model  $M_l$ , and the parameters  $\boldsymbol{\theta}$  so that

$$[\mathbf{y} | M_l] = \int [\mathbf{y} | \boldsymbol{\theta}, M_l][\boldsymbol{\theta}]d\boldsymbol{\theta}. \quad (10)$$

Note that this (Eq. 10) is the same expression typically appearing in the denominator of Bayes rule (Eq. 8).

#### Bayes factors

Assuming that we can find the posterior distribution for the quantity of interest  $[g | \mathbf{y}, M_l]$  for all models in  $\mathcal{M}$ , we need only compute the posterior model weights to find the averaged posterior distribution (Eq. 7). As it happens, solving the integral in the marginal data distribution (Eq. 10) is often nontrivial, which is why most Bayesian studies use MCMC to avoid calculating it directly. The sum in the denominator of the posterior model probability (Eq. 9) can also become intractable as the number of models  $L$  grows. Thus, despite its attractiveness and rigor, the challenge with BMA is in its implementation.

Consider the ratio of posterior probabilities for two models, say  $M_l$  and  $M_{\bar{l}}$ . Using a bit of algebra it is easy to show that the ratio (i.e., the posterior odds) is

$$\begin{aligned} \frac{P(M_l | \mathbf{y})}{P(M_{\bar{l}} | \mathbf{y})} &= \frac{[\mathbf{y} | M_l]P(M_l)}{\sum_{j=1}^L [\mathbf{y} | M_j]P(M_j)} \bigg/ \frac{[\mathbf{y} | M_{\bar{l}}]P(M_{\bar{l}})}{\sum_{j=1}^L [\mathbf{y} | M_j]P(M_j)} \\ &= \frac{[\mathbf{y} | M_l]P(M_l)}{[\mathbf{y} | M_{\bar{l}}]P(M_{\bar{l}})} = B_{l,\bar{l}} \frac{P(M_l)}{P(M_{\bar{l}})} \end{aligned} \quad (11)$$

which, after the data  $\mathbf{y}$  have been observed, can be written as a constant multiplier of the ratio of prior model probabilities (i.e., the prior odds). The multiplier  $B_{l,\bar{l}}$  in Eq. 11 is known as the Bayes factor and is only a function of the marginal data distributions from each model (Kass and Raftery 1995). Thus, the posterior evidence in favor of one model over another is found by updating the prior evidence with the data. Similar to the various rules of thumb for comparing models using

information criteria, there have been several suggested rules of thumb in the literature for Bayes factors (e.g.,  $B_{l,\bar{l}} > 10$  implies strong evidence in favor of model  $M_l$  over model  $M_{\bar{l}}$  according to Jeffreys (1961)).

The utility of the marginal data distribution for model averaging becomes clear because the posterior probability of any model  $M_l$

$$P(M_l | \mathbf{y}) = \frac{B_{l,\bar{l}}P(M_l)}{\sum_{j=1}^L B_{j,\bar{l}}P(M_j)} \quad (12)$$

is obtained by dividing the numerator and denominator in the posterior model probability (Eq. 9) by  $[\mathbf{y} | M_{\bar{l}}]$  (Link and Barker 2006). Thus, if we have the marginal data distributions  $[\mathbf{y} | M_l]$  for all models being considered, then we have the Bayes factors  $B_{l,\bar{l}}$ , and if we have the Bayes factors we can compute the exact Bayesian model weights for performing model averaging. Various methods exist for calculating the necessary quantities in Bayesian model averaging (e.g., Congdon 2006), some of which we will describe in what follows (*Statistical regularization and information criteria: Traditional regularizer: The penalty: Bayesian information criterion and Model-based model selection: Reversible-jump MCMC*). Finally, we note that one must be cautious in Bayesian model averaging when improper priors (i.e., prior distributions that do not integrate to 1) are used for parameters, as the Bayes factors are undefined in those settings (Spiegelhalter and Smith 1982).

#### Willow Tit occupancy: BMA

Royle and Dorazio (2008) describe a data set involving occupancy sampling of Swiss breeding birds as part of the Swiss Survey of Common Breeding Birds (collected by the Swiss Monitoring Haufige Brutvogel, and originally provided by Hans Schmid and Marc Kery). Thanks to Royle and Dorazio (2008), these data have become a standard textbook example used to demonstrate Bayesian occupancy models. We use a subset of data consisting of the first 200 quadrats throughout Switzerland where surveys were conducted for up to three sampling occasions. We focus on the same species considered by Royle and Dorazio (2008), the Willow Tit (*Parus montanus*), a relatively common passerine in Europe that resembles the Chickadee of North America in appearance (see Plate 1). Royle and Dorazio (2008) analyzed a binary form of the data at each site and occasion (i.e., detected/non-detected) along with covariate information on elevation and forest cover (which we standardize to have mean zero and standard deviation equal to one). Further details concerning data collection methods for this study are described by Kery and Schmid (2004).

Existing life history information concerning the environmental niche of the willow tit suggests that forest cover and elevation are important features. To demonstrate Bayesian model averaging (as well as the methods that follow) applied to the occupancy model,

TABLE 2. Prior and posterior model probabilities for Willow Tit occupancy.

Model	Covariates	$P(M_l)$	$P(M_l \mathbf{y})$
$M_1$	NULL	0.25	0.00
$M_2$	ELEV	0.25	0.52
$M_3$	FOR	0.25	0.00
$M_4$	ELEV + FOR	0.25	0.48

Note:  $M_1$ – $M_4$  are models,  $P(M_l)$  is the prior model probability,  $P(M_l|\mathbf{y})$  is the posterior probability of model  $M_l$ , and  $\mathbf{y}$  denotes a data set, typically composed of response variables. NULL, ELEV, FOR, and ELEV + FOR are the intercept-only, elevation, forest, and elevation plus forest models, respectively.

we constructed a set of four distinct candidate models to learn about the niche preferences of this species. Each occupancy model contains a homogeneous detection probability and an occupancy probability that (1) is homogeneous, containing only an intercept (i.e.,  $\beta_0$ ; the null model,  $M_1$ ), (2) contains an intercept and elevation as a covariate ( $M_2$ ), (3) contains an intercept and forest as a covariate ( $M_3$ ), and (4) contains an intercept and both elevation and forest as covariates ( $M_4$ ).

Assuming that we seek to use within-sample data to combine models, we can utilize Bayesian model averaging to obtain improved inference concerning the niche preferences of Willow Tit in Switzerland. Using the computational approaches described in *Model-based model selection* (i.e., reversible-jump MCMC), we calculated the posterior model probabilities for the four models ( $M_1$ – $M_4$ ; Table 2). Assuming equal prior probabilities for this example (i.e.,  $M_l = 1/4$  for  $l = 1, \dots, 4$ ), we find that the two models containing the elevation covariate dominate the model averaged inference with posterior model probabilities of  $P(M_2|\mathbf{y}) = 0.52$  and  $P(M_4|\mathbf{y}) = 0.48$ . Given our equal prior model probabilities, the Bayes factor for model  $M_2$  over  $M_4$  is computed as  $P(M_2|\mathbf{y})/P(M_4|\mathbf{y}) = 1.08$ .

We demonstrate the differences between posterior means for coefficients among all models considered in Table 3 as well as the model averaged posterior means. Notice that the BMA posterior mean for the elevation coefficient falls between the values resulting from the two models containing that covariate (i.e.,  $M_2$  and  $M_4$ ), while the BMA posterior mean for the forest coefficient shrinks toward zero. This shrinkage of  $\beta_1$  is caused by the very small posterior model probability for  $M_2$  (i.e., the model with only forest as a covariate), thus down weighting the estimate resulting from that model because it carries little weight in the Bayesian model average.

Following the line of reasoning provided by Madigan and Raftery (1994) it is common to consider BMA for only the two models containing the elevation covariate because the others have negligible posterior model probabilities. Thus, if one desired BMA inference based on the Occam’s window principle (i.e., considering only models carrying substantial weight in the averaging),

one would rerun the analysis using only the two top models in this scenario. We return to Bayesian model averaging in *Model-based model selection*, describing various approaches for computation.

## MODEL VALIDATION

In this section, assume again that we are considering a set of models  $\mathcal{M}$ . But now suppose we are interested in evaluating each model’s performance relative to some predefined characteristic. Predictive ability is by far the most commonly sought model characteristic in the literature on model selection and thus we highlight it here. Alternatively, other methods have been developed for selection based on estimation inference (i.e., inference that seeks to improve our understanding of model parameters rather than predictions [Bondell and Reich 2012]).

### Out-of-sample validation

If we are interested in prediction as our main characteristic of model utility, then it is sensible to evaluate the model in terms of real predictive ability; that is, we seek a model whose predictions are close to out-of-sample (oos) data (with closeness measured using a score function). Out-of-sample data are observations that are not used to fit the model but that we can use to compare with model predictions. In the machine learning literature, out-of-sample data are often referred to as “validation” data, whereas within-sample data are commonly referred to as “training” data (Hastie et al. 2009).

The essential idea in out-of-sample validation is that two data sets are collected; one to fit (or train) the model ( $\mathbf{y}$ ) and one to validate the model ( $\mathbf{y}_{\text{oos}}$ ). A large out-of-sample data set will provide the best information about the predictive performance of a model, but is obviously more intensive to collect. Thus, some trade-off between within-sample and out-of-sample data set size is necessary. For large single data sets such as those derived from web searches or financial data it is common to split the data set into two pieces, one for training and another for validation. If the original data set is large enough, the resulting decrease in inferential power due to splitting it up is negligible. In historical ecological studies it was less common to have such large

TABLE 3. Willow Tit occupancy posterior means for  $p$ ,  $\beta_0$ ,  $\beta_1$ , and  $\beta_2$  across all models and using Bayesian model averaging (BMA).

Parameter	$M_1$	$M_2$	$M_3$	$M_4$	BMA
$p$	0.26	0.26	0.26	0.26	0.26
$\beta_0$	0.17	0.38	0.89	0.29	0.32
$\beta_1$	0.00	1.95	0.00	1.80	1.85
$\beta_2$	0.00	0.00	1.79	0.39	0.18

Note: The parameters are  $p$ , detection probability;  $\beta_0$ , the intercept;  $\beta_1$ , the slope due to elevation; and  $\beta_2$ , the slope due to forest.

data sets, at least in terms of response variables. However, with remote sensing and newer automated data collection methods such as global positioning system (GPS) telemetry devices, large ecological data sets are more common than ever. Thus, out-of-sample validation methods are becoming more realistic for ecological analyses.

Out-of-sample validation relies on the ability to compute a similarity statistic or scoring rule to obtain a measure of closeness between our out-of-sample data  $\mathbf{y}_{\text{OOS}}$  and the predictions  $\hat{\mathbf{y}}_{\text{OOS}}$  (e.g., Bernardo 1979, Czado et al. 2009, Gneiting and Raftery 2007, Gneiting 2011). One of the most commonly used scoring rules is the mean squared prediction error (MSPE)

$$\text{MSPE} = \sum_{i=1}^{n_{\text{OOS}}} \frac{(y_{i,\text{OOS}} - \hat{y}_{i,\text{OOS}})^2}{n_{\text{OOS}}} \quad (13)$$

or its square root (RMSPE). The prediction,  $\hat{y}_{i,\text{OOS}}$  in MSPE, is obtained without using the out-of-sample observation  $y_{i,\text{OOS}}$ . The out-of-sample validation procedure can be applied independently for each model in a discrete set of models  $\mathcal{M}$  and the predictive scores (e.g.,  $\text{RMSPE}_i$  for model  $M_i$ ) can be compared to assess which model is best overall at prediction or how the models rank in terms of predictive ability.

The MSPE is a popular scoring rule because it has important properties when used with certain models. In general, Bernardo and Smith (1994) recommend logarithmic scoring rules that are both “local” and “proper.” In essence, these scoring rule characteristics guarantee that the predictive score adheres to the chosen model and data (Vehtari and Ojanen 2012, Gelman et al. 2014b). We describe a more general approach for scoring models based on out-of-sample data in what follows.

The practice of evaluating models based only on point estimates of parameters or predictions does not naturally incorporate our uncertainty pertaining to those quantities. One of the primary advantages of Bayesian inference is the ability to account for various sources of uncertainty, thus we now describe a method for model validation that appropriately accommodates uncertainty. In doing so, it is critical to recall how prediction works from the Bayesian perspective. In general, data that have not been observed are considered to be random quantities, thus we treat them like all other random quantities in the Bayesian setting and seek their posterior distribution. The posterior distribution for predictions is called the “posterior predictive distribution” and can be found using the integral

$$[\mathbf{y}_{\text{OOS}} | \mathbf{y}] = \int [\mathbf{y}_{\text{OOS}} | \mathbf{y}, \boldsymbol{\theta}] [\boldsymbol{\theta} | \mathbf{y}] d\boldsymbol{\theta}. \quad (14)$$

One option for the point prediction itself ( $\hat{\mathbf{y}}_{\text{OOS}}$ ) could be the posterior predictive mean, which technically requires another integral. That is

$$\hat{\mathbf{y}}_{\text{OOS}} = E(\mathbf{y}_{\text{OOS}} | \mathbf{y}) = \int \int \mathbf{y}_{\text{OOS}} [\mathbf{y}_{\text{OOS}} | \mathbf{y}, \boldsymbol{\theta}] [\boldsymbol{\theta} | \mathbf{y}] d\boldsymbol{\theta} d\mathbf{y}_{\text{OOS}} \quad (15)$$

which can be easily approximated as long as the out of sample data  $\mathbf{y}_{\text{OOS}}$  can be sampled from the distribution  $[\mathbf{y}_{\text{OOS}} | \mathbf{y}, \boldsymbol{\theta}]$  within an MCMC algorithm. If this condition is met, one can use composition sampling (Tanner 1996) and Monte Carlo integration to approximate the point prediction by

$$\hat{\mathbf{y}}_{\text{OOS}} \approx \frac{\sum_{t=1}^T \mathbf{y}_{\text{OOS}}^{(t)}}{T} \quad (16)$$

where  $\mathbf{y}_{\text{OOS}}^{(t)}$  is the  $t$ th MCMC sample (out of  $T$  total MCMC samples) of the predicted out-of-sample data. That is, we draw  $\mathbf{y}_{\text{OOS}}^{(t)}$  as a sample from  $[\mathbf{y}_{\text{OOS}} | \mathbf{y}, \boldsymbol{\theta}^{(t)}]$  at every MCMC iteration  $t$  for  $t = 1, \dots, T$  and then average them.

The procedure we have just described provides a way to obtain Bayesian point predictions, but it does not directly accommodate uncertainty pertaining to a score function. As it turns out, the log predictive density  $\log [\mathbf{y}_{\text{OOS}} | \mathbf{y}]$  is a local and proper scoring function that is appropriate for Bayesian model validation (Gelman et al. 2014b). In the situation where we have actual out-of-sample data  $\mathbf{y}_{\text{OOS}}$ , then we could just compute

$$\log \left( \frac{\sum_{t=1}^T [\mathbf{y}_{\text{OOS}} | \mathbf{y}, \boldsymbol{\theta}^{(t)}]}{T} \right) \quad (17)$$

using MCMC samples  $\boldsymbol{\theta}^{(t)}$ , as a Monte Carlo integral representation of the score function

$$\log [\mathbf{y}_{\text{OOS}} | \mathbf{y}] = \log \int [\mathbf{y}_{\text{OOS}} | \mathbf{y}, \boldsymbol{\theta}] [\boldsymbol{\theta} | \mathbf{y}] d\boldsymbol{\theta}. \quad (18)$$

This score can then be used to rank all models in the set  $\mathcal{M}$  and find the one that yields the best predictions. Out-of-sample validation is almost as efficient as simply fitting the individual models because it only requires the additional calculation of  $[\mathbf{y}_{\text{OOS}} | \mathbf{y}, \boldsymbol{\theta}^{(t)}]$  on each MCMC iteration, which is a low-order operation. Thus, for large ecological data sets, the out-of-sample validation approach is a very reasonable way to find good predictive models. However, as the out-of-sample size reduces, this validation procedure becomes less stable and thus more sensitive to the set of out-of-sample data.

#### Cross-validation

The concept of cross-validation was developed as a way to increase the stability of validation based on out-of-sample data for smaller sample sizes. Cross-validation is similar to out-of-sample validation in that we exclude a subset of the data ( $\mathbf{y}_k$ ) from the fitting procedure so that the model is unaware of it, and then compute the score based on the excluded data. The problem with choosing a single subset of the data to leave out is that you can only assess predictive ability for



those measurements. Thus, it is common to leave out all of the data, but only in small subsets sequentially.

$K$ -fold cross-validation involves grouping the data evenly (or approximately even) into  $K$  groups and then using each set of left out data  $\mathbf{y}_k$  to compare with the model predictions based on the remaining data ( $\mathbf{y}_{-k}$ ). We then iterate through all groups of data  $\mathbf{y}_k$  for  $k = 1, \dots, K$  and compute component scores which are summed to yield the full cross-validation score for the whole data set

$$\sum_{k=1}^K \log \left( \frac{\sum_{t=1}^T [\mathbf{y}_k | \mathbf{y}_{-k}, \boldsymbol{\theta}^{(t)}]}{T} \right). \quad (19)$$

In the case where  $K = n$  ( $n$  is the sample size), the procedure is often referred to as leave-one-out cross-validation. Leave-one-out cross-validation may be preferable when the sample size is small and there are few observations to use as training data, though the resulting estimate of prediction error becomes less stable as  $K \rightarrow n$ .

In general, the major disadvantage of  $K$ -fold cross-validation for Bayesian models is that we are required to refit each statistical model  $K$  times to obtain the complete set of out-of-sample predictions. Acquiring  $K \times L$  individual model fits may be reasonable for simple models, but for more complicated models that take longer to fit, a  $K$ -fold increase in required computing time may not be reasonable. However, despite these challenges, when true predictive ability is the main criterion of interest, cross-validation is still very appealing for model comparison. In fact, it underlies several parsimony-based model comparison methods.

#### Conditional predictive ordinates

To improve computational tractability for large data and model sets, one could consider the posterior predictive distribution for within-sample data. That is, instead of cross-validation, simply compute the aforementioned predictive score based on the predictive distributions of the data  $[y_i | \mathbf{y}]$  for  $i = 1, \dots, n$ . The problem with this approach is that the predictive performance of the model will be overestimated because the data are used twice (i.e., once for model fitting and another time for model validation). The overestimation of predictive performance is referred to as “optimism” in the statistics literature and we return to this concept in Section 4.

As a potential remedy, consider the leave-one-out predictive distribution for each observation in a data set

$$[y_i | \mathbf{y}_{-i}] = \int [y_i | \boldsymbol{\theta} | \mathbf{y}_{-i}] d\boldsymbol{\theta}. \quad (20)$$

This quantity in Eq. 20 is referred to as the conditional predictive ordinate (CPO<sub>*i*</sub>; Geisser 1993) and represents the probability (or density) of the observation  $y_i$  when the model is fit without that observation. Thus, large CPO<sub>*i*</sub> values correspond to

very likely observations under the current model, whereas small CPO<sub>*i*</sub> indicates outliers and/or high-leverage observations (Pettit 1990). In principle, the computation of CPO would require a true cross-validation involving an  $n$ -fold iterative model fitting scheme. Fortunately, CPO can be approximated easily within an MCMC algorithm for model fitting as the harmonic mean of the predictive distributions evaluated at the MCMC values for the parameters  $\boldsymbol{\theta}$

$$\text{CPO}_i \approx \frac{T}{\sum_{t=1}^T [y_i | \boldsymbol{\theta}^{(t)}]^{-1}} \quad (21)$$

where  $t = 1, \dots, T$  represent the MCMC iterations. A summary statistic of these individual CPO values, such as  $-\sum_i \log(\text{CPO}_i)$ , then provides an overall measure of predictive performance. Notice the similarity in expressions for the sum of the logged CPO values and the log predictive score (Eq. 19) described in the previous section. In terms of appropriateness for model selection, the CPO involves a harmonic mean, which yields a numerically unstable estimator in practice, but software can often be constructed to flag problematic cases (Held et al. 2010).

#### Willow Tit occupancy: model validation

Suppose that we are now interested in comparing the four occupancy models we introduced in *Model validation* in terms of their predictive ability. We do not have an auxiliary source of out-of-sample data to use for model validation, but we can employ Bayesian cross-validation and also compute the  $-\sum_i \log(\text{CPO}_i)$  statistic based on Eq. 21 to compare the information about predictive ability using each of these methods.

We used 10-fold Bayesian cross-validation (i.e.,  $K = 10$ ) due to the moderate sample size and computed the scoring function discussed in Eq. 19 as

$$-2 \sum_{k=1}^{10} \log \left( \frac{\sum_{t=1}^T \text{Binom}(\mathbf{y}_k | \mathbf{J}_k, p^{(t)} \mathbf{z}_k^{(t)})}{T} \right) \quad (22)$$

where,  $p^{(t)}$  and  $\mathbf{z}_k^{(t)}$  are MCMC samples arising from model fits not including observations  $\mathbf{y}_k$  and the  $-2$  is multiplied merely for convenience (so that small scores are better and to compare with other model selection criteria later). Thus, the inner sum in Eq. 22 is over the MCMC iterations from a single fold of the validation procedure and the outer sum is over the  $K$  folds. We obtained 160 000 MCMC iterations to fit each model (in each fold), discarding the first 16 000 as burn-in. To illustrate the computational gains achieved using contemporary parallel programming methods we performed the cross-validation using both non-parallel and parallel algorithms. The non-parallel algorithm (i.e., a single loop over the  $K$  folds) required approximately 1 hour, whereas the parallel algorithm required over an order of magnitude less computing time at approximately 5.7 minutes. Similarly, it required 1.4 minutes to compute the CPO statistics in parallel (but 5.7 minutes in

TABLE 4. Willow Tit occupancy results for cross-validation and CPO.

Model	Covariates	Cross-validation score	$-\sum_i \log(\text{CPO}_i)$
$M_1$	NULL	552.4	240.2
$M_2$	ELEV	478.4	220.0
$M_3$	FOR	526.9	246.2
$M_4$	ELEV + FOR	478.8	220.4

sequence). All computation was performed on a desktop workstation with two 2.93 GHz 6-Core processors and 32 GB of RAM; we note that new laptops have individual processors that are substantially faster, but parallel computing is still more efficient on the desktop we used with its many cores. All MCMC algorithms were coded natively in R (R Core Team 2013) and the R package snowfall (Knaus 2013) was used for parallel computing.

In Table 4, we can see that the Bayesian cross-validation score generally agrees with CPO in that the two models with elevation as a covariate (i.e.,  $M_2$  and  $M_4$ ) out-perform the null model ( $M_1$ ) and model with only an intercept and forest as a covariate ( $M_3$ ; note also that lower scores are better). The null model performs the worst based on the cross-validation score, while the two models with elevation are nearly equivalent in terms of prediction. CPO indicates that the null model may be slightly better at prediction than the model with only forest as a covariate (i.e.,  $M_3$ ), however, given that cross-validation evaluates predictive performance based on out-of-sample data, we might be skeptical of these CPO results for the worst performing models. This potential discrepancy between cross-validation and CPO is part of the sacrifice we make when computation time is limited.

#### STATISTICAL REGULARIZATION AND INFORMATION CRITERIA

The assessment of a set of models in terms of their predictive ability has been a central theme in the development of information criteria. However, information criteria involve specific approaches to model selection that fall under the much broader umbrella of statistical regularization. This concept of regularization, though used on a daily basis in ecology, does not appear to be widely recognized. However, regularization reveals numerous theoretical and practical connections among model selection and multimodel inference paradigms. Specifically, regularization links Bayesian and non-Bayesian approaches to model selection and here we describe how this linkage occurs. We begin by presenting the basic regularization concept, showing how it has been used traditionally in the non-Bayesian context. We then describe how regularization is inherently Bayesian and highlight a few explicitly Bayesian approaches for doing it (e.g., the Bayesian Lasso).

The term “regularization” refers to the use of an external regulator that constrains the results of an optimization problem (note that the term “regulator” is borrowed here from physics but is not commonly used

in statistics, though it is perhaps more intuitive). In statistical terminology, the optimization problem could be a likelihood that needs maximizing or a posterior distribution that needs exploring (perhaps via MCMC). In the broader decision theoretic context, we might refer to a negative log-likelihood more generically as a loss function; that is, a function that expresses the “loss” incurred by inadequately estimating parameters of interest. In certain cases, the loss function may have too much freedom to be useful for inference and thus an external constraint can help make it useful.

In placing this concept of regularization in a formal statistical framework for decision making, or parameter estimation, consider the generic expression

$$\mathcal{L}(\mathbf{y}, \boldsymbol{\theta}) + r(\boldsymbol{\theta}, \gamma). \quad (23)$$

$\mathcal{L}(\mathbf{y}, \boldsymbol{\theta})$  represents the loss, a function of both knowns ( $\mathbf{y}$ ) and unknowns ( $\boldsymbol{\theta}$ ) and, though it is related, should not to be confused with a likelihood (which we label  $\mathbf{y} | \boldsymbol{\theta}$ ). The function  $r(\boldsymbol{\theta}, \gamma)$  in Eq. 23 represents the regulator or constraint on the unknowns  $\boldsymbol{\theta}$ . The regulator function  $r$  may also depend on some other variables  $\gamma$  that may or may not be related to the loss function or its components. There are other ways to express the loss and regulator relationship, but the expression in Eq. 23 is perhaps the most common. Statistical inference can now be obtained by minimizing the joint function (Eq. 23) with respect to  $\boldsymbol{\theta}$ , and perhaps  $\gamma$ , if not already known. The primary advantage of regularization is that it can yield improved inference, often reducing the variance of estimates and increasing the accuracy of predictions. Though not often discussed in the ecological literature, this concept of regularization is quite common in many areas of statistics and machine learning (Hastie et al. 2009). As we will see in the next sections, regularization also underlies the dominant model selection approaches used in ecology and has direct ties with Bayesian statistics.

#### *Traditional regulator: the penalty*

To make the concept of regularization more concrete, we place it in the context of classical non-Bayesian regression modeling. That is, consider the linear model

$$y_i \sim \mathbf{N}(\beta_0 + \mathbf{x}'_i \boldsymbol{\beta}, \sigma^2) \quad (24)$$

where  $\mathbf{N}$  denotes the normal distribution for  $i = 1, \dots, n$ , where the “unknowns” are the regression coefficients  $\beta_0$  and  $\boldsymbol{\beta}$ . For now, assume the error variance  $\sigma^2$  is known, but note that it need not be in general. If our goal is to find estimates of  $\beta_0$  and  $\boldsymbol{\beta}$ , then the loss function for this optimization problem is proportional to the negative log-likelihood  $\mathcal{L}(\mathbf{y}, \beta_0, \boldsymbol{\beta}) = \sum_{i=1}^n (y_i - \beta_0 - \mathbf{x}'_i \boldsymbol{\beta})^2$ . Now consider the regulator function  $\gamma_1 \sum_{j=1}^p |\beta_j|^2$ , called the “penalty” in the statistical literature, such that the optimization problem from Eq. 23 becomes

$$\sum_{i=1}^n (y_i - \beta_0 - \mathbf{x}'_i \boldsymbol{\beta})^2 + \gamma_1 \sum_{j=1}^p |\beta_j|^2 \quad (25)$$

where  $p$  corresponds to the dimension of  $\boldsymbol{\beta}$  (i.e., the number of covariates in the model),  $\gamma_1$  is often referred to as the penalization or bandwidth parameter (in the statistics literature,  $\lambda$  is often used instead of  $\gamma_2$ ; we avoid the  $\lambda$  notation here to reduce any confusion with the leading eigenvalue of a Leslie matrix in demographic modeling), and the exponent  $\gamma_2$  is the chosen degree of the “norm.” Note that the penalty is commonly written using norm notation, that is,  $\|\boldsymbol{\beta}\|_{\gamma_2} \equiv \sum_{j=1}^p |\beta_j|^{\gamma_2}$  (referred to as the  $L_{\gamma_2}$  norm for a specific value of  $\gamma_2$ ). The parameters  $\gamma_1$  and  $\gamma_2$  control the amount and type of regularization that occurs in the estimation problem. Although the parameters  $\gamma_1$  and  $\gamma_2$  are sometimes chosen only implicitly, based on adherence to a particular philosophical underpinning, there seems to be greater variety in the rationale and practical choices for  $\gamma_1$  than for  $\gamma_2$ . We discuss commonly used choices for  $\gamma_2$  next.

*Ridge regression.*—So-called “ridge regression” is a direct application of above optimization problem in Eq. 25 where the parameter  $\gamma_2 = 2$  is used in the penalty term. In this case, we seek to minimize

$$\sum_{i=1}^n (y_i - \beta_0 - \mathbf{x}'_i \boldsymbol{\beta})^2 + \gamma_1 \sum_{j=1}^p \beta_j^2 \quad (26)$$

with respect to the regression coefficients  $\beta_0$  and  $\boldsymbol{\beta}$  given a certain value for the penalty parameter  $\gamma_1$ . If  $\gamma_1 = 0$ , then the negative log-likelihood is not penalized and the resulting estimated coefficients will be the maximum likelihood estimates (MLEs). However, as  $\gamma_1$  increases, it will “shrink” the estimated coefficients  $\boldsymbol{\beta}$  toward zero when Eq. 26 is minimized as a trade-off between maximizing the likelihood and meeting the constraint. This is why regularization methods in the maximum likelihood setting are commonly referred to as “penalized” or “shrinkage” methods. The shrinkage of  $\boldsymbol{\beta}$  can be incredibly useful in parameter estimation and prediction.

In parameter estimation, shrinkage induces an increasing bias in  $\hat{\boldsymbol{\beta}}$  with increasing  $\gamma_1$  but simultaneously reduces the variance of  $\hat{\boldsymbol{\beta}}$ . Thus, in ridge regression, we accept a small amount of bias in our estimation of  $\boldsymbol{\beta}$  in return for a potentially large reduction in variance. The reduction in variance of  $\hat{\boldsymbol{\beta}}$  also decreases prediction error, providing improved prediction accuracy. More complex models provide an excellent fit to within-sample data but are poor predictors of out-of-sample data. Shrinking model parameters toward zero reduces effective model complexity thereby improving our ability to predict out-of-sample data.

These features of ridge regression are undoubtedly desirable, but may overshadow one of the most useful aspects of the regularization: alleviation of the effect of multicollinearity in the covariates (e.g., Graham 2003). When columns of our “design matrix”  $\mathbf{X}$  are correlated with each other, the associated coefficients  $\boldsymbol{\beta}$  have to compete for the overall effect on the response variables

and this competition causes the coefficient estimates  $\hat{\boldsymbol{\beta}}$  to offset each other, forcing some to be very large (positive) and some very small (negative). In cases where significant multicollinearity exists, the penalty term in the optimization problem will shrink these exaggerated parameter estimates back to reasonable values. Thus, in ridge regression, we can use the “full” model including all the variables in  $\mathbf{X}$  at once, regardless of how much they are correlated with each other. The alternative approach is to construct a finite model set where no single model contains any two covariates that are correlated beyond a certain threshold (e.g., correlation coefficient  $\rho = 0.6$ , as advocated by Burnham and Anderson [2002]). This latter approach is a type of discrete regularization, rather than a continuous one such as ridge regression.

There are a few practical considerations in the proper application of regularization methods for regression models. First, notice that we have separated the intercept  $\beta_0$  from the rest of the regression coefficients  $\boldsymbol{\beta}$  in Eq. 25. We isolate  $\beta_0$  because we do not wish to shrink the general mean of the regression model to zero, rather, only the coefficients that interact with covariates. Second, it is advisable to standardize the covariates in  $\mathbf{X}$  prior to analysis (i.e., subtract the mean and divide by the standard deviation). This standardization of covariates allows us to use a single penalty parameter  $\gamma_1$  rather than one for each coefficient  $\beta_j$  so that they do not need to be shrunk differentially. The third consideration is the choice of  $\gamma_1$ , which we discuss in the next section.

*Lasso: Least absolute shrinkage and selection operator.*—Continuing with the linear regression example (Eq. 25) used in the previous section, now consider a different regulator function where we set  $\gamma_2 = 1$  such that

$$\sum_{i=1}^n (y_i - \beta_0 - \mathbf{x}'_i \boldsymbol{\beta})^2 + \gamma_1 \sum_{j=1}^p |\beta_j|. \quad (27)$$

This new penalty term ( $\gamma_1 \sum_{j=1}^p |\beta_j|$ ) is commonly referred as the “Lasso” or  $L_1$ -norm penalty and induces a markedly different constraint on the optimization problem. The acronym Lasso stands for least absolute shrinkage and selection operator (Tibshirani 1996) because the use of an  $L_1$ -norm penalty implies a sum of absolute coefficient values. While the  $L_2$ -norm penalty in ridge regression shrinks  $\boldsymbol{\beta}$  toward zero nonlinearly (with increasing  $\gamma_1$ ), the  $L_1$ -norm Lasso penalty shrinks the coefficients linearly in such a way that they eventually can equal zero exactly in the optimization. Thus, Lasso drops covariates from the model by setting their coefficients to zero. This absolute variable selection concept seems quite familiar to many ecologists who learned about model selection from a traditional perspective. This heuristic familiarity has made the Lasso approach very popular (Dahlgren 2010).

To summarize, we have now seen that both  $\gamma_1$  and  $\gamma_2$  in Eq. 25 play important roles in statistical regularization. Given that  $\gamma_1$  controls the amount of shrinkage

induced, it acts as a type of *scale* parameter, while  $\gamma_2$  controls the form of the shrinkage and could be thought of as a *shape* parameter. For now, we suspect that the choice of  $\gamma_2$  is more a result of personal preference based on desired inference, but what about  $\gamma_1$ ? How should we choose the amount of shrinkage?

Heuristically, we seek inference concerning model parameters that is based on a balance between model fit and predictive ability. Thus, we could treat  $\gamma_1$  as we do any other model parameter and estimate it simultaneously with the others. The problems with this approach are manifold, but relate to the same basic concept: within-sample data vs. out-of-sample data. Even if there is enough information in the data to actually estimate an “extra” model parameter, the fact that within-sample data are being used to learn about  $\gamma_1$  limits its utility as a regulator. Recall from our discussion of cross-validation, that there are trade-offs in using the same set of data to both fit and validate (i.e., select) models. The primary trade-off is that predictive performance can only truly be assessed using out-of-sample data. Thus, it seems most reasonable to estimate model parameters based on within-sample data and choose regulator parameters based on out-of-sample data.

A strategy employed in many machine learning studies is to optimize the regularized loss function (Eq. 23) given the within-sample data  $\mathbf{y}$  for the first term and use an iterative cross-validation approach to choose  $\gamma_1$  based on predictive ability of out-of-sample data. In practice, a strategy for the regression model would involve first optimizing Eq. 25 using  $\gamma_1 = 0$  assigning a cross-validation score, and then incrementally increasing  $\gamma_1$  over a range of values yielding a set of predictive scores. Given a sufficiently fine range of values for  $\gamma_1$ , we would then choose the regularized model yielding the best predictive score. In the case of ridge regression, our inference would consist of a full set of coefficient estimates  $\hat{\boldsymbol{\beta}}$  that are properly shrunk to provide the best predictions of out-of-sample data. For Lasso, we would obtain a subset of non-zero coefficient estimates that have been shrunk according to the  $L_1$ -norm penalty, and the remaining coefficients would be zero (i.e., no longer in the final model). In either case, we will obtain a justifiably parsimonious model that is better at prediction than the unpenalized full model. Another advantage is that we did not have to do prior variable elimination based on highly collinear covariate pairs.

Despite the many advantages to classical regularization, there are also several disadvantages. Aside from the somewhat ad hoc and subjective feel of the procedure, these methods are based on optimization and they yield point estimates for the model parameters of interest, but learning about the uncertainty of  $\hat{\boldsymbol{\beta}}$  is not necessarily trivial or even possible in some cases. Finally, because we may want to rely on out-of-sample data to choose appropriate regulator parameters ( $\boldsymbol{\gamma}$ ), this can dramatically increase the computational requirements of cross-validation-based regularization.

*Akaike's information criterion.*—Continuing in a non-Bayesian context, we now explain how information criteria fit into the regularization concept. Statistical regularization is appealing for the reasons discussed in the previous section, but for many ecologists, the increased computational burden and need to select regulator parameters can be daunting. Enter the information criterion approach to statistical regularization. The general idea behind information criteria is that we choose a scoring function a priori that will be used to “score” each of the models based on the balance of fit using the within-sample data and parsimony (or overall predictive ability; Gneiting 2011). Not surprisingly, most commonly used information criteria take the same form as the previously introduced regularization expression (Eq. 23). For example, in the linear regression class of models, Akaike's information criterion (AIC) takes the form of Eq. 25 with regulator parameters  $\gamma_1 = 2$  and  $\gamma_2 = 0$  such that the penalty is  $2\sum_{j=1}^p |\beta_j|^0 = 2p$ . The  $L_0$ -norm used in AIC implies that the shrinkage is only based on the number of parameters rather than the parameter values themselves. This implication is useful because each model in the model set can be fit independently and then post hoc scored using AIC (lower AIC implying better predictive ability of the model). However, we must be careful to avoid inducing obvious bias in the estimates by choosing a model set such that no single model contains correlated covariates because the penalty cannot provide feedback to the estimation of the parameters themselves.

AIC provides the same regularization as leave-one-out cross-validation under certain conditions (Stone 1977). We find this a very appealing result on first glance because it could dramatically reduce the computational burden in finding a good predictive model. However, upon closer inspection, we find that the result only holds in linear Gaussian settings (i.e., regression models with additive normal errors) and under the assumption that the “true” model is in the model set being considered. This latter assumption (i.e., truth in the model set) seems to conflict with one of the main advantages of AIC extolled by proponents. Still, empirically, AIC seems to perform well in situations where it can be used (Hastie et al. 2009). For Bayesians, AIC (being a function of maximum likelihood estimates) does not appear to have a clear Bayesian interpretation, at least outside of a few contrived situations (as we discuss later in *Bayesian regulator: the prior*).

The use of an information criterion like AIC requires a compromise: We trade the continuous aspects of model selection using more general regulators (e.g., ridge regression, Lasso) for the reduction in computational burden achieved by avoiding cross-validation.

*Bayesian information criterion.*—The so-called Bayesian information criterion (BIC; Schwarz 1978) arises from a different motivation than does AIC and many other regularization methods. AIC is an information



criterion that seeks to provide a measure of predictive ability, whereas BIC is distinctly concerned with multi-model inference (Link and Barker 2006, Gelman et al. 2014b).

Recall the marginal data distribution  $[y | M_l]$  for model  $M_l$  from *Model averaging* (Eq. 10). The marginal data distribution is critical for computing Bayes factors and model probabilities in the Bayesian paradigm. In a maximum likelihood setting, if we consider the loss function to be  $-2 \log[y | \hat{\theta}]$ , as is assumed with AIC, then we can approximate the marginal data distribution using a Laplace approximation (Ripley 1996) such that for model  $M_l$

$$\text{BIC} = -2\log[y | \hat{\theta}, M_l] + \log(n)p \approx -2\log[y | M_l] \quad (28)$$

where  $\log(n)$  is the natural logarithm of the sample size (or dimension of  $y$ ) and  $p$  is the number of “free” parameters, as before. Note that, for the linear regression model (Eq. 24), this definition of BIC still retains the general regularization form of (Eq. 25), but with regulator parameters  $\gamma_1 = \log(n)$  and  $\gamma_2 = 0$ .

The utility of BIC in multimodel inference arises when we exponentiate negative one-half times the BIC (Eq. 28); normalizing this quantity over all models in the model set  $\mathcal{M}$  provides an approximation to the Bayesian model weights (Eq. 9) described previously. Unfortunately, this approximation only holds when equal prior model weights (i.e.,  $P(M_l) = 1/L$  for  $l = 1, \dots, L$ ) are assumed. Furthermore, because of its reliance on maximum likelihood parameter estimates, BIC does not appear to be inherently Bayesian (despite its name). Finally, BIC can only be used to approximate posterior model probabilities when the Bayes factors are well defined, which is not the case if improper priors are used in the models.

From a classical perspective, there is no clear choice, nor consensus, among statisticians, between AIC and BIC for model selection purposes (Hastie et al. 2009). Each form of automatic regulator has advantages and disadvantages. For example, BIC can be shown to be a consistent model selector (i.e., the oracle property). That is, when the “true” model is in the model set and the data set is sufficiently large, BIC will select the true model, while AIC will select models that are too large in general. On the other hand, for smaller sample sizes, BIC may indicate models that are too parsimonious because  $\log(n) > 2$  implies more shrinkage from BIC than AIC. Furthermore, BIC is motivated from a model averaging rather than prediction perspective, and thus it may be more justified for approximating Bayesian model weights than for model selection.

#### *Bayesian regulator: the prior*

The previous section describes regularization from a classical perspective, where we penalize a statistical optimization problem in such a way that it yields a better predictive model. As we hinted at earlier, the

fact that the classical regularization approach seems to “work” is encouraging, but its lack of formality brings up a set of new questions (e.g., What type of regulator function to use? How much shrinkage is too much?). Furthermore, on the surface, the classical regularization methods do not appear to be able to accommodate uncertainty about the parameters or regulator function. For ecologists using Bayesian models, what is the analog to regularization in the Bayesian setting?

*Natural Bayesian shrinkage.*—The analog to regularization in the Bayesian setting is simply the Bayesian model itself! To see this, consider the linear regression example (Eq. 24) used in the previous section, but now, we specify priors for the unknown model parameters  $\beta$  such that the model itself is specified as

$$y_i \sim N(\beta_0 + \mathbf{x}'_i \beta, \sigma^2)$$

$$\beta \sim N(\mu, \sigma_\beta^2 \mathbf{I}) \quad (29)$$

where, for illustrative purposes, we assume the intercept  $\beta_0$  and variance parameter  $\sigma^2$  are fixed and known for now. The posterior distribution for  $\beta$  is then easily shown to be

$$[\beta | y] \propto [y | \beta] [\beta]$$

$$\propto \prod_{i=1}^n N(y_i | \beta_0 + \mathbf{x}'_i \beta, \sigma^2) \prod_{j=1}^p N(\beta_j | \mu_j, \sigma_\beta^2)$$

$$\propto \exp\left(-\frac{1}{2} \frac{\sum_{i=1}^n (y_i - \beta_0 - \mathbf{x}'_i \beta)^2}{\sigma^2}\right) \exp\left(-\frac{1}{2} \frac{\sum_{j=1}^p (\beta_j - \mu_j)^2}{\sigma_\beta^2}\right)$$

$$\propto \exp\left(-\frac{1}{2} \left(\frac{\sum_{i=1}^n (y_i - \beta_0 - \mathbf{x}'_i \beta)^2}{\sigma^2} + \frac{\sum_{j=1}^p (\beta_j - \mu_j)^2}{\sigma_\beta^2}\right)\right)$$

$$\propto \exp\left(-\frac{1}{2\sigma^2} \left(\sum_{i=1}^n (y_i - \beta_0 - \mathbf{x}'_i \beta)^2 + \frac{\sigma^2}{\sigma_\beta^2} \sum_{j=1}^p (\beta_j - \mu_j)^2\right)\right). \quad (30)$$

If we let  $\mu_j = 0$  for all  $j = 1, \dots, p$ , and reparameterize the ratio of variances such that  $\gamma_1 = \sigma^2/\sigma_\beta^2$  in the last expression of Eq. 30, then we arrive at the exact same regularization expression used in ridge regression (Eq. 26) in the inner parentheses of our posterior distribution for  $\beta$  (Eq. 30). Thus, by reducing our prior variance for the regression coefficients, we increase the effective regulator parameter  $\gamma_1$  and induce the same sort of shrinkage on  $\beta$  as in ridge regression, but in a formal Bayesian probability framework. In fact, one could say that we are always doing a form of regularization in Bayesian statistics because the prior acts as the regulator. Given that the Bayesian posterior provides a rigorous framework for regularization, it could be argued that other classical forms of regularization are inherently Bayesian, or at least Bayesian in spirit.

Regardless of the interpretation of the regulator, as a non-Bayesian penalty or as a Bayesian prior, we can enjoy the same benefits of regularization from either perspective. However, the Bayesian perspective makes it clear that we are constraining the model parameters with “prior” information such that it assists us in finding a better predictive model. We are often taught that the Bayesian prior should either be chosen objectively as to minimize the influence on the posterior, or retrospectively, to best represent existing prior knowledge about the parameters. However, the only rule for specifying prior information in a Bayesian model is to not use the within-sample data to choose the prior. The reason for this rule is that it maintains the acyclicity in the Bayesian “graph.” Bayesian models are often referred to as directed acyclic graphs because of their conditional specifications such that the data depend on the parameters and the parameters depend on either other parameters or fixed quantities. The acyclic nature of the Bayesian graph guarantees that we can use valid probability statements to learn about the unknown quantities. Interestingly, this rule of “don’t use the data twice” is commonly broken, and the model is referred to as empirical Bayesian in that setting. Empirical Bayesian methods seem to perform well, as does classical regularization, but have much weaker theoretical foundations than fully Bayesian methods. It seems clear that to fit a rigorous Bayesian model we should not use the within-sample data in the likelihood and the prior, but there is no such rule about the use of out-of-sample data to inform the prior. Thus, we could think of the three ways to specify valid priors as (1) objectively, (2) retrospectively, and (3) prospectively. The term “prospective” in this sense implies the use of future data, perhaps collected at the same time as the within-sample data but not used until after (rather than before) the likelihood is specified. This third approach to specifying priors opens up the door for Bayesian cross-validation.

For example, the Bayesian cross-validation procedure for regularization of the regression model might proceed as follows: Specify the model as in Eq. 29, fit it for each of the  $K$  sets of hold-out data using a vague prior for  $\beta$  with mean zero and obtain a predictive score as described in *Model validation: Cross-validation*. Choose an incrementally smaller prior variance  $\sigma_\beta^2$  and repeat the model fitting and cross-validation scoring process. Continue this procedure, using smaller and smaller prior variances until an optimal predictive model is identified (typically via a small score function). Finally, fit the optimal predictive Bayesian regression model using the full data set to obtain desired inference.

The problem arises in the last step of this cross-validation procedure. Once we use the prior (i.e., penalty or regulator) that has been informed by an aggregate of hold-out data, we technically cannot put all of the hold-out data back into the model to fit one last time for final inference in a fully Bayesian paradigm. In this case, the options are to use the data twice in this way and accept

that the procedure is empirical Bayesian, or use two completely separate data sets, one for training ( $\mathbf{y}$ ) and another for validating ( $\mathbf{y}_{\text{oss}}$ ). Of course, the second option is not always preferable when analyzing data that have already been collected, but in larger data sets or when setting up new studies, collecting two independent data sets for two different purposes allows for fully rigorous Bayesian inference and model selection.

*Bayesian lasso.*—The previous section illustrates how the standard Bayesian regression model with a Gaussian prior on the coefficients provides a natural mechanism to perform statistical regularization similar to ridge regression, but how can we manipulate the regulator function? The answer is simple in the regression case: We only need to find a prior with the same form as the desired regulator function. For example, to construct a Bayesian regularization that has a penalty similar to the Lasso penalty, we need only find a prior containing an  $L_1$ -norm on the parameters. In this case, the Laplace distribution contains the  $L_1$ -norm that will impose a Lasso penalty as a prior. That is, consider the same regression data model, but with a new prior for  $\beta$  such that

$$y_i \sim N(\beta_0 + \mathbf{x}'_i \beta, \sigma^2)$$

$$\beta_j \sim \text{Laplace}(\mu = 0, \sigma_\beta^2) \propto \exp\left(-\frac{|\beta_j|}{\sqrt{\sigma_\beta^2}}\right) \quad (31)$$

for  $j = 1, \dots, p$  where  $\beta_j$  are independent a priori. Park and Casella (2008) propose a similar prior for  $\beta$ , as well as more standard priors for  $\beta_0$  and  $\sigma^2$  and dub it “the Bayesian Lasso.” In fact, they go a step further and carefully specify a prior for a transformation of the regulator parameter that enables them to construct a fully conjugate MCMC algorithm for fitting the model. Unlike in a Metropolis-Hastings MCMC algorithm, the resulting Gibbs sampler requires no tuning of any parameters (Kyung et al. 2010). Thus, it is nearly as computationally efficient to fit the Bayesian Lasso regression model as it is the standard Bayesian regression model. Of course, Bayesian cross-validation could also be used in this scenario and would likely yield better out-of-sample predictive performance, but would also require substantially more computational effort.

Finally, after seeing the connection between Bayesian priors and regulator functions, one might wonder what sort of prior yields an AIC penalty? Following the same approach described in the Bayesian Lasso (Eq. 31), it appears that the implicit AIC prior for each coefficient is  $[\beta_j] \propto \exp(-|\beta_j|^0)$ , such that the joint prior distribution for  $\beta$  is  $[\beta] \propto \exp(-p)$ .

#### *Willow Tit occupancy: Bayesian regularization*

In applying Bayesian regularization to the Willow Tit occupancy model, we first remind the reader that the model already contains a natural regularization mechanism: the prior for  $\beta$ . Recall the process component of

the hierarchical occupancy model from Eq. 3

$$v_i \sim N(\beta_0 + \mathbf{x}'_i \boldsymbol{\beta}, 1) \quad (32)$$

and prior from Eq. 6

$$\boldsymbol{\beta} \sim N(\boldsymbol{\mu}_\beta, \sigma_\beta^2 \mathbf{I}). \quad (33)$$

Notice that if we standardize the covariates to have mean zero and variance one then we can reasonably set the prior mean  $\boldsymbol{\mu}_\beta = \mathbf{0}$ . In this case, the full-conditional distribution for  $\boldsymbol{\beta}$  becomes

$$[\boldsymbol{\beta} | \cdot] \propto \exp\left(-\frac{1}{2} \left( \sum_{i=1}^n (v_i - \beta_0 - \mathbf{x}'_i \boldsymbol{\beta})^2 + \frac{1}{\sigma_\beta^2} \sum_{j=1}^p \beta_j^2 \right)\right) \quad (34)$$

as was demonstrated for the regression model (Eq. 30). Thus, this full-conditional distribution for  $\boldsymbol{\beta}$  has the same form as the general regularization expression (Eq. 25) and the hyperparameter  $\sigma_\beta^2$  serves as the regulator or shrinkage parameter, where  $\gamma_1 = 1/\sigma_\beta^2$ . In other words, the smaller we make the prior variance, the stronger the penalty in the regularization. The strategy is to explore the space of  $\sigma_\beta^2$  for an optimal value that provides the best predictive model according to the score function of choice. To find the optimal penalty, we can explore the space of  $\sigma_\beta^2$  using a grid search (i.e., try a range of  $n_\beta$  total values for  $\sigma_\beta^2$ ) and compare scores based on cross-validation. This cross-validation approach requires  $K \times n_\beta$  separate model fits, resulting in a potentially unreasonable amount of required computational time. For example, a 10-fold cross-validation, at 1.4 minutes per model fit and  $n_\beta = 24$  dimensional grid search would require 5.6 hours to implement. However, using 24 processors in parallel, the required time could be reduced to under an hour on a high-performance desktop workstation. The three easy ways to reduce computation time are to (1) use more processors (e.g., a high-performance computing facility), (2) decrease the number of folds in the cross-validation (e.g., an  $n$ -fold cross-validation for the above example would require almost 5 days in sequence, but only a few hours in parallel) and (3) use a lower-resolution grid search. The latter will require fewer model fits on the same machine, but will reduce the accuracy of the optimization.

We wouldn't expect Bayesian regularization to dramatically increase predictive ability for the simple willow tit occupancy model because the two covariates (elevation and forest) are relatively uncorrelated (i.e., correlation  $\approx 0.12$ ) and the sample size ( $n = 200$ ) is large relative to the number of unknown parameters. However, to demonstrate the regularization approach, we use the full model for the Willow Tit data with one intercept and two regression coefficients associated with the occupancy probability ( $M_4$ ). We then perform a grid search over 24 values for  $\sigma_\beta^2$ , implying a prior

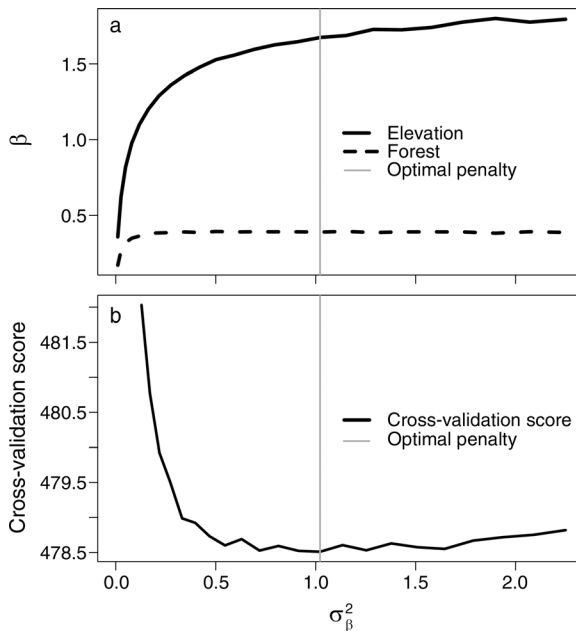


FIG. 3. Willow Tit occupancy: Bayesian regularization. (a) Shrinkage trajectories for the posterior mean of model parameter  $\boldsymbol{\beta}$  ( $y$ -axis) plotted against prior variance for  $\boldsymbol{\beta}$  ( $x$ -axis). Parameter estimates yielding the best predictive model based on the two covariates occur at the vertical gray line. Note that the correlation between elevation and forest is 0.12. (b) The cross-validation score ( $y$ -axis) presented in Eq. 22 plotted against prior variance for  $\boldsymbol{\beta}$  ( $x$ -axis). The optimal score (i.e., smallest; score = 478.5) for prediction occurs at the vertical gray line (i.e., minimum score occurs at  $\sigma_\beta^2 = 1.02$ ).

that ranges from precise ( $\sigma_\beta^2 = 0.01$ ) to vague ( $\sigma_\beta^2 = 2.25$ ).

We used the log posterior predictive score for 10-fold cross-validation introduced earlier (22). The complete 10-fold cross-validation at each value of  $\sigma_\beta^2$ , with model fits based on 160 000 MCMC iterations (discarding 16 000 as burn-in), took approximately 24 minutes with parallel computing.

We found that the optimal prior variance for prediction occurs at  $\sigma_\beta^2 = 1.02$ ; this is less than half of the variance we would typically use in a vague prior scenario for the occupancy model. In Fig. 3 we see the posterior means for  $\boldsymbol{\beta}$  taper toward zero as  $\sigma_\beta^2$  decreases. At the optimal level of regularization, the predictive score was 478.4, yielding a model that predicts as well as  $M_2$  (the elevation only model) but uses both covariates. Notice also that the cross-validation score function increases more sharply away from the optimum as  $\sigma_\beta^2$  decreases toward zero. This effect indicates that the null model (i.e., occurring at  $\sigma_\beta^2 = 0$ ) performs substantially worse than the full model (i.e., occurring at  $\sigma_\beta^2 = 2.25$ ), a result similar to that found in the former cross-validation of the discrete model set (Table 4).

*Deviance information criterion*

We have seen that a natural framework for regularization in the Bayesian context already exists and can be used in conjunction with out-of-sample data to help select an appropriate penalty. However, the classical information criteria were developed, at least in part, to alleviate the need for cross-validation and seem to perform quite well in many settings. Is there a Bayesian equivalent?

Spiegelhalter et al. (2002) proposed the deviance information criterion (DIC), which has a similar form as other information criteria, in that it contains a loss function plus a penalty or regulator function. The loss function is chosen to be the deviance

$$D(\boldsymbol{\theta}) = -2\log[\mathbf{y} | \boldsymbol{\theta}] \quad (35)$$

as in most other information criteria, but in order to be similar to AIC or BIC the penalty needs to incorporate the number of free parameters as a measure of model complexity. Recall that, even in the simplest Bayesian models, most parameters are constrained in some way by their priors. Furthermore, in hierarchical Bayesian models, we may have numerous latent state variables that are technically unknown but are also highly constrained by both the likelihood and prior. Thus, one crucial issue in the development of a truly Bayesian criterion is the specification of an “effective” number of parameters, say  $p_D$ . A further complication is that maximum likelihood point estimates are used to compute AIC and BIC, but this concept of maximum likelihood is only meaningful under certain situations in the Bayesian context. Thus, we can use a Bayesian point estimate, the posterior mean, in lieu of the MLE in DIC

$$\text{DIC} = -2\log[\mathbf{y} | E(\boldsymbol{\theta} | \mathbf{y})] + 2p_D = \hat{D} + 2p_D \quad (36)$$

where  $E(\boldsymbol{\theta} | \mathbf{y})$  corresponds to the posterior expectation of  $\boldsymbol{\theta}$  and the deviance evaluated at the posterior mean for  $\boldsymbol{\theta}$  is commonly written as  $\hat{D}$ .

To arrive at a measure of model complexity, Spiegelhalter et al. (2002) consider the difference in the deviance calculated two different ways: posterior mean deviance and deviance computed at the posterior mean of the parameters. That is, the effective number of parameters was originally defined as

$$p_D = \bar{D} - \hat{D} \quad (37)$$

such that the posterior mean deviance is

$$\bar{D} = E_{\boldsymbol{\theta} | \mathbf{y}}(-2\log[\mathbf{y} | \boldsymbol{\theta}]) = \int -2\log[\mathbf{y} | \boldsymbol{\theta}][\boldsymbol{\theta} | \mathbf{y}]d\boldsymbol{\theta}. \quad (38)$$

In the case of linear regression, with vague priors on the regression coefficients, the effective number of parameters  $p_D$  approaches the number of coefficients  $p$ . Thus, the popularity of DIC has been a result of its similarity to AIC, its simplicity, and its ease of calculation using MCMC samples. There are only two quantities that need to be computed for DIC: The

deviance evaluated at the posterior mean of the parameter set  $\hat{D}$ , which is as trivial as the deviance calculation in AIC, and the posterior mean deviance, which can be embedded into an MCMC algorithm with one or two lines of code.

For many Bayesian models (which we describe in the next section), DIC can be used for ranking models and finding those that should predict better than others, just as AIC would. DIC addresses the issue of model complexity and in many cases yields results quite similar to AIC. A common question is whether DIC can be used for Bayesian model averaging? That is, if one follows the AIC-based guidance of Burnham and Anderson (2002), and calculates  $w_j = e^{-\Delta\text{DIC}_j/2} / \sum_l e^{-\Delta\text{DIC}_l/2}$ , where  $\Delta\text{DIC}_j$  represents the difference of DIC for model  $j$  and the minimum DIC across all models in the model set, do these weights  $w_j$  approximate posterior model probabilities? Despite the fact that this approach is used occasionally, the answer has not been justified in the literature. Link and Barker (2006) make a strong case for the use of BIC to approximate posterior model probabilities and perform a small set of empirical comparisons between AIC, BIC, and DIC model weighting schemes, but the theoretical foundation for Bayesian model averaging using DIC is much weaker.

*Modified DIC.*—Despite its convenience, DIC has several limitations, notable among them are the potential for poorly estimating model complexity ( $p_D$ ), inappropriateness with mixture models, and the lack of a direct connection with predictive ability. We elaborate on some of these issues with conventional DIC before discussing some attractive alternatives.

There have been many alternative specifications for the effective number of parameters  $p_D$  (Eq. 37), which is sometimes referred to as model complexity, or degrees of freedom, in the statistical literature. For example, Plummer (2002) suggests that a more appropriate measure of model complexity can be computed by averaging

$$\log \left( \frac{[\tilde{\mathbf{y}}^{(1,t)} | \boldsymbol{\theta}^{(1,t)}]}{[\tilde{\mathbf{y}}^{(2,t)} | \boldsymbol{\theta}^{(2,t)}]} \right) \quad (39)$$

over all MCMC samples (i.e.,  $t = 1, \dots, T$ ), where  $\tilde{\mathbf{y}}^{(1,t)}$  and  $\tilde{\mathbf{y}}^{(2,t)}$  are two independent posterior predictive realizations of the data arising from two different chains (for  $\boldsymbol{\theta}^{(1,t)}$  and  $\boldsymbol{\theta}^{(2,t)}$ ) based on separate model fits. This version of model complexity (Eq. 39) arises as an estimate of the expected Kullback-Leibler divergence between predictive distributions at two values for  $\boldsymbol{\theta}$  (Plummer 2002). Unfortunately, Plummer (2008) later indicates that the average of Eq. 39 may only be an appropriate penalty when the sample size is very large (i.e.,  $n \rightarrow \infty$ ). Plummer (2008) also recommends an alternative estimator for model complexity with better properties, but its calculation requires  $n$  separate model fits, which puts it on par with cross-validation, thus reducing the appeal of DIC in terms of computational





PLATE 1. Willow tit (*Parus montanus*) observed in the Swiss Alps. Note that an alternate genus name for the willow tit is *Poecile*. Photo credit: copyright© Marcel Burkhardt, used with permission.

efficiency. Overall, it appears that DIC (Eq. 36) is most appropriate as a model selection criterion in linear models with independent data (conditional on  $\theta$ ) where the  $p_D$  is much smaller than  $n$ . Thus, DIC is good for comparing Bayesian versions of the same classes of models that AIC is good for comparing.

Several others have suggested that DIC is not appropriate for model selection with mixture models or missing data models (e.g., Spiegelhalter et al. 2002, Celeux et al. 2006, Plummer 2008). Zero-inflated models comprise the largest and most heavily used class of models in wildlife ecology (i.e., capture–recapture and occupancy models) and are a form of mixture model (Martin et al. 2005). The original version of DIC is thus not suitable for comparing zero-inflated models. Celeux et al. (2006) provide several suggestions that could be used as an alternative to the standard DIC for mixture models, but ultimately they do not recommend any of them as a gold standard. However, one of these modified versions of DIC was also discussed earlier by Richardson (2002) and lacked a theoretical justification until recently (Watanabe 2010). Celeux et al. (2006) numbered this information criterion  $DIC_3$ , and we discuss it next.

#### *Watanabe-Akaike information criterion*

Aside from the aforementioned caveats, DIC is a useful information criterion in the parametric Bayesian modeling context when prediction is of primary impor-

tance. However, DIC does not best represent the actual Bayesian predictive procedure. To arrive at predictions, the Bayesian approach is to find and summarize the posterior predictive distribution (Eq. 14). In computing DIC (Eq. 36), the posterior predictive distribution is not needed. This seems to be a mismatch between the type of inference desired and the tool used to obtain it.

Along the same lines of reasoning we used in the previous section on out-of-sample validation, for Bayesian model comparison based on predictive ability, we should seek a statistic that considers the log posterior predictive distribution for new data  $\tilde{y}$

$$\log[\tilde{y} | \mathbf{y}] = \log \int [\tilde{y} | \theta][\theta | \mathbf{y}] d\theta. \quad (40)$$

The quantity in Eq. 40 is stochastic because  $\tilde{y}$  is assumed to be unknown (but not so in true out-of-sample validation scenarios; hence the change in notation from  $y_{\text{OOS}}$  to  $\tilde{y}$ ), therefore a common technique in the development of most information criteria is to then consider the mean of Eq. 40 over  $\tilde{y}$

$$E_{\tilde{y}}(\log[\tilde{y} | \mathbf{y}]) = \int \log \int [\tilde{y} | \theta][\theta | \mathbf{y}] d\theta d\tilde{y} \quad (41)$$

which is impossible to compute directly because the true distribution of the new data  $[\tilde{y}]$  is unknown. Thus, in finding an estimator of mean log posterior predictive

score, Richardson (2002), Celeux et al. (2006), and Watanabe (2010) propose the log point-wise predictive score

$$\log \prod_{i=1}^n [y_i | \mathbf{y}] = \sum_{i=1}^n \log \int [y_i | \boldsymbol{\theta}] [\boldsymbol{\theta} | \mathbf{y}] d\boldsymbol{\theta} \quad (42)$$

where Monte Carlo integration can be used to compute the integral (Gelman et al. 2014b). There are two issues with the score in Eq. 42: (1) the product representation of the posterior predictive distribution implies that the data are independent (conditioned on  $\boldsymbol{\theta}$ ) and (2) it relies completely on the observed data  $\mathbf{y}$  rather than the new data  $\tilde{\mathbf{y}}$ . The first issue suggests that the score should not be used with models containing dependence in the data (e.g., spatial and time series models). The latter issue implies that Eq. 42 will be optimistic in its predictive score for a given model because the within-sample data are being used twice. As in DIC, the amount of optimism with this score (Eq. 42) can be expressed as the effective number of parameters  $p_D$  (Watanabe 2010). Thinking of the effective number of parameters  $p_D$  in this way is not intuitive because most ecologists have been trained to view the penalty in AIC as  $p$ , the actual number of parameters. In fact,  $p$  in that sense is really a measure of model complexity that arises naturally in the derivation of many information criteria. Thus, it is helpful to think of  $p_D$  as a measure of model complexity rather than strictly a count of the model parameters.

Gelman et al. (2014b) present two possible estimates for  $p_D$

$$p_{D,1} = 2 \sum_{i=1}^n \left( \log E_{\boldsymbol{\theta} | \mathbf{y}} [y_i | \boldsymbol{\theta}] - E_{\boldsymbol{\theta} | \mathbf{y}} (\log [y_i | \boldsymbol{\theta}]) \right) \quad (43)$$

and

$$p_{D,2} = \sum_{i=1}^n \text{var}_{\boldsymbol{\theta} | \mathbf{y}} (\log [y_i | \boldsymbol{\theta}]) \quad (44)$$

but prefers  $p_{D,2}$  for its relationship with leave-one-out cross-validation. As with DIC, we can use Monte Carlo integration to approximate  $p_{D,2}$  by computing the sum of the MCMC sample variances of  $\log [y_i | \boldsymbol{\theta}^{(t)}]$  (sample variance computed over  $t = 1, \dots, T$  MCMC samples) over the observations  $y_i$  for  $i = 1, \dots, n$ .

The Watanabe-Akaike information criterion can then be defined as  $-2$  times the log point-wise predictive score plus the estimated optimism

$$\text{WAIC} = -2 \sum_{i=1}^n \log \int [y_i | \boldsymbol{\theta}] [\boldsymbol{\theta} | \mathbf{y}] d\boldsymbol{\theta} + 2p_{D,2} \quad (45)$$

with both elements in the sum approximated using MCMC samples at no extra computational cost beyond that required for calculating DIC (Watanabe 2013). The addition of the estimated optimism in Eq. 45 serves as a bias correction in estimating posterior predictive accuracy similar to that of AIC and DIC, even though we have not mentioned it until now. The term “optimism,”

which is often used in the statistical literature, is merely another word for regulator or penalty.

This new criterion enjoys many benefits. Among them are the fact that WAIC is based on the posterior predictive distribution and is fully Bayesian, but yields the same results as DIC in linear Gaussian models with uniform priors. Furthermore, unlike DIC, WAIC is valid in both hierarchical and mixture models (Watanabe 2013). Also, unlike DIC, the effective number of parameters calculated using  $p_{D,2}$  in Eq. 44 will always be positive. In  $p_{D,2}$ , a parameter gets counted as a 1 if all of the learning we gain about it comes from the likelihood. Conversely, a parameter counts as a zero in the calculation of  $p_{D,2}$  if the learning comes entirely from the prior. To figure out the correct proportion of each parameter to count, WAIC needs to use the data (like in DIC) to compute the optimism  $p_{D,2}$ . This is essential in the Bayesian context where we regularly use hierarchical structures with strong interdependencies and informative priors.

Overall, WAIC seems very appealing, however, the main disadvantage is substantial depending on the area of application: its calculation relies on an independence assumption of the data given the parameters. This assumption is regularly violated in spatial models where dependence among the data is one of the key features being modeled. Ando and Tsay (2010) provide a way to relax the independence assumption, but the resulting criterion requires numerous model fits, which eliminates one of the key practical benefits of WAIC (Gelman et al. 2014b).

#### Posterior predictive loss

In a similar spirit as that motivating WAIC, and in contrast with CPO, another approach to prediction-based model choice was presented by Laud and Ibrahim (1995) and later justified by Gelfand and Ghosh (1998). This approach, referred to as “posterior predictive loss,” considers prediction from a decision theoretic perspective. Understanding this approach requires a familiarity with statistical decision theory, which we describe briefly here, referring the interested reader to more comprehensive references (e.g., Berger 2006, Vehtari and Ojanen 2012) for further details.

Statistical decision theory provides a rigorous framework for the decision making process in the presence of data and uncertainty (Berger 2006). The phrase “decision-making process” is quite general, encompassing decisions like choices of alternatives for management, but also including a justification for parameter estimation and prediction. In fact, behind every statistical estimator lies a set of implicit or explicit decision theoretic assumptions. A formal decision theory exists in both the classical and Bayesian realms, though Berger (2006) makes a compelling case for the completeness of the Bayesian decision theory.

In essence, a Bayesian decision theory involves three main concepts: (1) a loss function, (2) an “action” or

decision, and (3) a posterior risk function. The loss function is a mathematical expression of the loss incurred if a certain decision is made and the posterior risk function is the loss averaged over the posterior distribution for the unknown quantities of interest. Thus, risk is a version of loss that has accounted for our uncertainty about the study system. The statistical literature refers to the decision minimizing the posterior risk as a “Bayes rule” (Lehmann and Casella 1998).

For example, suppose we are interested in estimating a parameter  $\theta$  given data  $\mathbf{y}$ . In the case of parameter estimation, the “decision” is actually just a point estimator of  $\theta$ . A point estimate  $\hat{\theta}$  that minimizes our risk seems desirable, thus the Bayes rule for point estimation is called a Bayes estimator. To find this Bayes estimator, we simply define a function  $\mathcal{L}(\mathbf{y}, \theta)$  that suitably represents the loss we incur for poorly estimating  $\theta$  and minimize its average with respect to the posterior distribution. The value for  $\theta$  that minimizes the posterior risk  $\hat{\theta}$  is the resulting Bayes estimator.

As it turns out, the Bayes estimator for squared error loss (i.e.,  $\mathcal{L}(\mathbf{y}, \theta) = (\theta - \hat{\theta})^2$ ) is the posterior mean of  $\theta$ , a result that we often use for inference without putting much thought into the rationale for why we use it. Different loss functions result in different estimators. For example, the absolute loss (i.e.,  $\mathcal{L}(\mathbf{y}, \theta) = |\theta - \hat{\theta}|$ ) results in the posterior median as the Bayes estimator and zero-one loss (i.e.,  $\mathcal{L}(\mathbf{y}, \theta) = 0$  or  $\mathcal{L}(\mathbf{y}, \theta) = 1$  if  $\theta = \hat{\theta}$  or  $\theta \neq \hat{\theta}$ , respectively) results in the posterior mode being the Bayes estimator.

Returning to the topic of model selection, Gelfand and Ghosh (1998) recommended a decision theoretic approach based on prediction rather than parameter estimation. In doing so, they proposed a loss function in terms of hypothetical replicates of the data  $\tilde{y}_i$  (i.e., unobserved new data) that is a sum of two components

$$\mathcal{L}(\tilde{y}_i, \hat{y}_i) + w\mathcal{L}(y_i, \hat{y}_i) \quad (46)$$

where  $\hat{y}_i$  represents a predictive realization for the unobserved new data point  $\tilde{y}_i$ , and  $y_i$  represents the observed within-sample data point. In the proposed loss function (Eq. 46), the  $w$  is constrained to be nonnegative and expresses the relative weight given to loss for the within-sample vs. new data at the same prediction  $\hat{y}_i$ .

Gelfand and Ghosh (1998) derived a posterior predictive risk by averaging their proposed loss function (46) over the posterior predictive distribution of  $\tilde{y}_i | \mathbf{y}$ . The resulting risk is then minimized with respect to the prediction  $\hat{y}_i$  and summed over all observations  $i = 1, \dots, n$  to yield the model selection criterion

$$D_w = \sum_{i=1}^n \min_{\hat{y}_i} \int (\mathcal{L}(\tilde{y}_i, \hat{y}_i) + w\mathcal{L}(y_i, \hat{y}_i)) [\tilde{y}_i | \mathbf{y}] d\tilde{y}_i \quad (47)$$

where we would seek to find a model with the smallest  $D_w$  out of a proposed set of models given a chosen loss

TABLE 5. Willow Tit occupancy results for WAIC, DIC, and  $D_{\infty, \text{sel}}$  (posterior predictive loss).

Model	Covariates	WAIC	DIC	$D_{\infty, \text{sel}}$
$M_1$	NULL	481.7	462.2	288.0
$M_2$	ELEV	440.2	432.2	270.8
$M_3$	FOR	492.4	483.8	305.2
$M_4$	ELEV + FOR	440.7	432.9	271.2

function  $L(\cdot)$  and weight  $w$ . In practice, it can be difficult to compute the necessary integrals in Eq. 47, thus a squared error loss (sel) function is commonly used, yielding the criterion

$$D_{w, \text{sel}} = \frac{w}{w+1} \sum_{i=1}^n (y_i - E(\tilde{y}_i | \mathbf{y}))^2 + \sum_{i=1}^n \text{Var}(\tilde{y}_i | \mathbf{y}). \quad (48)$$

Further, it is often assumed that the weight is very large ( $w \rightarrow \infty$ ) thus resulting in a  $D_{\infty, \text{sel}}$  criterion

$$D_{\infty, \text{sel}} = \sum_{i=1}^n (y_i - E(\tilde{y}_i | \mathbf{y}))^2 + \sum_{i=1}^n \text{Var}(\tilde{y}_i | \mathbf{y}). \quad (49)$$

Note the similarity of  $D_{\infty, \text{sel}}$  to the WAIC (Eq. 45) and DIC (Eq. 36, for large  $n$ ) in that they both contain two terms in a sum, the first being a goodness-of-fit measure and the second acting as a penalty or regulator. In this case, we can see that the penalty  $\sum_{i=1}^n \text{Var}(\tilde{y}_i | \mathbf{y})$  will increase in overfitted models where the prediction variance becomes larger with an increasing number of parameters.

For more general loss functions, such as deviance,  $D_w$  takes on a similar two component form, but the penalty is only guaranteed to be positive under certain constraints on the loss (i.e., convexity in  $y$ ) and the criterion may not be suitable for mixture models. Despite this caveat,  $D_w$  does appear to be appropriate for many classes of hierarchical models because it depends directly on the posterior predictive distribution rather than the likelihood and posterior mean of the parameters alone. Also, unlike WAIC, the general form of posterior predictive loss approach appears to be suitable for correlated data models (e.g., spatial and temporal models).

Even though the posterior predictive loss approach does not technically fall into the same category as the rest of the information criteria, the form of the general loss function proposed by Gelfand and Ghosh (1998) is similar enough to the regularization expression (Eq. 23), and equivalent to DIC and WAIC in certain settings, that we chose to describe it here rather than place it in its own section.

#### *Willow Tit occupancy: information criteria*

In a continued assessment of predictive performance for the occupancy model set using the Willow Tit data, we calculated WAIC, DIC, and  $D_{\infty, \text{sel}}$  for each of the 4 models previously considered (Table 5). To calculate WAIC for the occupancy model in this example, we used



MCMC samples to approximate the effective number of parameters

$p_{D,2} \approx$

$$\sum_{i=1}^n \frac{\sum_{t=1}^T \left( \log([y_i | J_i, p^{(t)} z_i^{(t)})] - \sum_{t=1}^T \log([y_i | J_i, p^{(t)} z_i^{(t)})] / T \right)^2}{T} \quad (50)$$

based on Eq. 44, where  $[y_i | J_i, p^{(t)} z_i^{(t)}]$  is the binomial probability mass function and the first term in WAIC (Eq. 45) is approximated as

$$-2 \sum_{i=1}^n \log \frac{\sum_{t=1}^T [y_i | J_i, p^{(t)} z_i^{(t)}]}{T}. \quad (51)$$

Recall that this expression (Eq. 51) has the same form as the cross-validation score (Eq. 22), but is based only on within-sample data.

For DIC, we used the traditional method for calculating the effective number of parameters (Eq. 37) and approximated  $\hat{D}$  and  $\hat{D}$  by

$$\hat{D} \approx \frac{\sum_{t=1}^T -2 \log[\mathbf{y} | \mathbf{J}, p^{(t)} \mathbf{z}^{(t)}]}{T} \quad (52)$$

$$\hat{D} \approx -2 \log[\mathbf{y} | \mathbf{J}, \hat{p} \hat{\mathbf{z}}] \quad (53)$$

where  $\hat{p}$  and  $\hat{\mathbf{z}}$  are the posterior means for detection probability and true latent occupancy status across all sites, and  $[\mathbf{y} | \mathbf{J}, p^{(t)} \mathbf{z}^{(t)}] = \prod_{i=1}^n [y_i | J_i, p^{(t)} z_i^{(t)}]$  is the likelihood based on the conditionally independent data for the willow tit occupancy model.

For the posterior predictive loss method, we calculated  $D_{\infty, \text{sel}}$  as in Eq. 49 based on the expectation and variance approximations

$$E(\bar{y}_i | \mathbf{y}) \approx \frac{\sum_{t=1}^T \bar{y}_i^{(t)}}{T} \quad (54)$$

$$\text{Var}(\bar{y}_i | \mathbf{y}) \approx \frac{\sum_{t=1}^T (\bar{y}_i^{(t)} - \sum_{t=1}^T \bar{y}_i^{(t)} / T)^2}{T} \quad (55)$$

where  $\bar{y}_i^{(t)} \sim [y_i | J_i, p^{(t)} z_i^{(t)}]$  is drawn on each MCMC iteration (for  $t = 1, \dots, T$ ) as a posterior predictive realization.

Of the three criteria considered in this example, recent statistical literature suggests that only WAIC is truly appropriate for the occupancy model (Gelman et al. 2014b). However, given that DIC is commonly used to compare Bayesian occupancy models, we provide a comparison here. Furthermore, the criterion based on posterior predictive loss ( $D_{\infty, \text{sel}}$ ) is not ideal for the occupancy model setting because the squared error loss function (Eq. 49) may not be best representative for the zero-inflated binomial data model. A different loss function could be chosen, but then a derivation would be required to find a computable approximation based on MCMC samples. Still, we felt that a comparison of the methods could illuminate potential empirical differ-

ences between the approaches. If this were a real application rather than a pedagogical example, we would have only computed WAIC for this model and data set. In terms of computational time, it only required 6.1 minutes to fit the models sequentially and obtain these metrics (using 160 000 MCMC iterations for each model fit with a burn-in period of 16 000 iterations).

All of these approaches (i.e., WAIC, DIC, and  $D_{\infty, \text{sel}}$ ) provide similar information in ranking the willow tit occupancy models by predictive ability based on within-sample data (Table 5). WAIC, DIC, and  $D_{\infty, \text{sel}}$  all suggest model  $M_3$ , the model containing only the forest covariate, as the worst predictive model, with the null model next ( $M_1$ ), and a virtual tie among the two models containing the elevation covariate (i.e.,  $M_2$  and  $M_4$ ). This latter result is in agreement with the earlier cross-validation and CPO model comparison.

#### MODEL-BASED MODEL SELECTION

To a certain extent, the regularization methods discussed in the previous section (especially the fully Bayesian Lasso) are model-based approaches to model selection. They are model based because they contain a formal mechanism that trades off model fit for model parsimony. We saw that the Bayesian model itself provides a natural model reduction mechanism via the prior. In contrast to this form of continuous shrinkage induced by a strong prior on the parameters, other methods have been developed in a similar spirit that explicitly augment the overall model structure with selection components whose job it is to switch on and off various effects in the full model (O'Hara and Sillanpaa 2009). The basic idea then is to build a model that contains all of the potential model components and then let the model decide which of them are helpful and which are not.

#### Indicator variable selection

For instructive purposes, consider again the basic linear regression model from Eq. 24

$$y_i \sim \mathbf{N}(\beta_0 + \mathbf{x}_i' \boldsymbol{\beta}, \sigma^2)$$

where, the parameter vector  $\boldsymbol{\beta} = (\beta_1, \dots, \beta_j, \dots, \beta_p)'$  contains the individual coefficients corresponding to the  $p$  predictor variables of interest. A modification of the original regression model has been proposed such that  $\beta_j = z_j \times \theta_j$  for  $j = 1, \dots, p$ , where each original parameter is written as a product of a binary indicator variable  $z_j$  and a regression coefficient  $\theta_j$  (e.g., George and McCulloch 1993, Carlin and Chib 1995, Kuo and Mallick 1998). In general, a prior would be specified for each  $(z_j, \theta_j)$  pair and the full Bayesian model could then be fit, yielding inference not only about the coefficients  $\beta_j$ , but also the selection indicators  $z_j$ . In this setting, if the posterior mean for a particular  $z_j$  is large (i.e., closer to one than zero) it would indicate that the  $j$ th covariate is important in the model; conversely,



when the posterior mean of  $z_j$  is close to zero it effectively removes the  $j$ th effect from the model thereby inducing a certain parsimony.

In implementing an indicator variable selection model, one would be tempted to use independent priors for  $z_j$  and  $\theta_j$ ; for example, we might specify

$$z_j \sim \text{Bern}(\phi)$$

$$\theta_j \sim \text{N}(0, \tau_j^2)$$

for all  $j = 1, \dots, p$ , assuming the covariates are standardized (where Bern stands for Bernoulli). However, an independent prior specification can cause computational problems if the prior for  $\theta_j$  is too vague (i.e., the prior variance,  $\tau_j^2$ , is large) because when  $z_j = 0$  in an MCMC algorithm,  $\theta_j$  will be sampled from its prior and the subsequent sampling of future  $z_j = 1$  will rarely occur since the  $\theta_j$  is likely to be far from the majority of posterior mass. Thus, to alleviate these computational problems, others (e.g., George and McCulloch 1993, Carlin and Chib 1995) have suggested joint priors for  $z_j$  and  $\theta_j$  that include explicit dependence between the indicators and coefficients.

In Gibbs variable selection, Carlin and Chib (1995) and Dellaportas et al. (1997) suggest decomposing the joint prior distribution  $[z_j, \theta_j] = [\theta_j | z_j][z_j]$ . In this joint prior specification, the Bernoulli prior for  $z_j$  is retained, but the prior for  $\theta_j$  conditional on  $z_j$  is written as

$$\theta_j | z_j \sim z_j \text{N}(0, \tau^2) + (1 - z_j) \text{N}(\mu_{\text{tune}}, \sigma_{\text{tune}}^2) \quad (56)$$

which has the form of a mixture distribution and is often referred to as a “slab and spike” prior (Miller 2002). The Gibbs variable selection procedure then involves choosing the tuning parameters  $\mu_{\text{tune}}$  and  $\sigma_{\text{tune}}^2$  such that  $\text{N}(\mu_{\text{tune}}, \sigma_{\text{tune}}^2)$  is near the posterior so that the MCMC algorithm exhibits better mixing. Surprisingly, the seemingly informative prior (Eq. 56) does not actually influence the posterior for  $\beta_j$ , but rather only influences the behavior of the MCMC algorithm (Carlin and Chib 1995).

In a similar model-based approach called “stochastic search variable selection,” George and McCulloch (1993) proposed a joint prior for  $z_j$  and  $\theta_j$ . However, unlike in the Gibbs variable selection, this alternative prior does influence the posterior and can be written as

$$\theta_j | z_j \sim z_j \text{N}(0, c\tau^2) + (1 - z_j) \text{N}(0, \tau^2). \quad (57)$$

In stochastic search variable selection, both  $c$  and  $\tau^2$  are tuned such that  $\tau^2$  is quite small, providing an effective spike at zero while  $c\tau^2$  is larger, creating a slab around zero. The slab then provides the prior for  $\theta_j$  when the variable  $\beta_j$  is in the model (i.e., when  $z_j = 1$ ). Both Gibbs and stochastic search variable selection methods require tuning to ensure well-mixed MCMC algorithms, but both can be useful for model-based model selection.

### Reversible-jump MCMC

A related model-based approach to model selection is referred to as reversible-jump Markov chain Monte Carlo (RJMC MC; Green 1995). Normally, we reserve the names of computational approaches for algorithms only, not statistical models; however, in this case, the method really describes a model, but we retain the label RJMC MC for convention. In describing the RJMC MC approach, first recall the model set  $\{M_1, \dots, M_I, \dots, M_L\}$  described earlier in *Model averaging: The utility of the marginal data distribution*. Now suppose that each of the models contain their own corresponding parameters  $\theta_l$ . Note that the lengths, say  $p_l$ , of these parameter vectors  $\theta_l$  may vary. In RJMC MC, we treat the model index  $l$  as a random quantity to be modeled along with the set of all possible parameters  $\theta$ . Or alternatively, we treat the number of parameters  $p_l$  as a random quantity and specify a model for it. Under certain assumptions, the posterior distribution of interest then is

$$[\theta, l | \mathbf{y}] \propto [y | \theta, l][\theta_l | l][l] \quad (58)$$

where  $[\theta_l | l]$  is the prior distribution for the parameters in model  $M_l$  and  $[l]$  is the prior distribution for model  $M_l$  itself. The beauty of this specification is that it places multimodel inference directly in a fully Bayesian context.

The use of MCMC to implement this model (Eq. 58) involves the usual steps: specify initial values for unknowns and then cycle through the unknowns, updating each one sequentially. The complication arises when sampling the model index  $l$ , and hence its associated parameters  $\theta_l$ , because the model dimension changes depending on which model is sampled. Thus, care must be taken to account for the potentially different model dimension when accepting a Metropolis-Hastings proposal for the parameters in an MCMC algorithm. The term “reversible” derives from the fact that certain properties of the Metropolis-Hastings update must be retained to arrive at a valid posterior distribution (Green 1995, Godsill 2001). Specifically, if we leave one model space with a particular dimension for another of a different dimension, we need to ensure that we can revert back to the former dimension later in the Markov chain. Thus, a modified version of the Metropolis-Hastings ratio can be constructed for certain models that corrects for the transdimensional nature of the algorithm.

RJMC MC approaches have become a popular option for computing Bayes factors and Bayesian model probabilities (e.g., Johnson and Hoeting 2011). When prior model probabilities are assumed to be equal, the Bayes factor ( $B_{l,i}$ ) can be computed simply by calculating the quotient of summed number of visits to each model ( $M_l$  and  $M_i$ ) in the RJMC MC algorithm (Hastie and Green 2012).

Due to its model-based form, RJMC MC is an appealing method for Bayesian multimodel inference

but can be tricky or impossible to implement for complicated models. To that end, Barker and Link (2013) described a method that provides RJMCMC results using a post hoc approach that only requires one to fit the  $L$  individual models and then post-process the resulting MCMC samples using a second MCMC algorithm in the form of a Gibbs sampler. We describe this approach and apply it to the willow tit data next.

In the big picture, Godsill (2001) and O'Hara and Sillanpaa (2009) show that the RJMCMC and indicator variable selection approaches are related. The key difference is that the auxiliary variables  $z_j$  are effectively moving the model between dimensions by switching on and off model components. In doing so, Gibbs and stochastic search variable selection side-step the trans-dimensional complication altogether.

#### *Willow Tit occupancy: RJMCMC*

We presented results pertaining to Bayesian model averaging earlier in *Model averaging*. To compute those Bayesian model averaging quantities, we used the RJMCMC approach described by Barker and Link (2013), which we briefly summarize here. One advantage of the Barker and Link (2013) approach is that the individual models can be fit separately and then recombined subsequently with a secondary MCMC algorithm to obtain posterior model probabilities. After the initial set of four occupancy models were fit individually (requiring only 5.7 minutes in sequence), the following secondary algorithm was constructed to iteratively sample the model and associated parameters:

- 1) Set MCMC iteration index to  $t = 1$ .
- 2) Choose initial model  $M_t^{(t)}$ . In our case we used  $M_t^{(1)} = M_4$ , the full model.
- 3) Select  $p_t^{(t)}$ ,  $\beta_{0,t}^{(t)}$ , and  $\beta_t^{(t)}$  from the former MCMC output for model  $M_t^{(t)}$ .
- 4) If there are remaining parameters from the full model not obtained in step 3 (i.e., for models  $M_1$ ,  $M_2$ , and  $M_3$ ) then sample those from a known distribution (the form of which is arbitrary according to Barker and Link 2013). We used a standard normal distribution to sample remaining parameters,  $N(0, 1)$ .
- 5) Order the parameter values from steps 3 and 4 and combine to form  $\theta$ . For example, if  $M_t^{(t)} = M_2$ , then  $\theta \equiv (p_t^{(t)}, \beta_{0,t}^{(t)}, \beta_{1,t}^{(t)}, u_2^{(t)})'$ , where  $u_2^{(t)} \sim N(0, 1)$ .
- 6) Compute the full-conditional model probability

$$P(M_l | \cdot) = \frac{[\mathbf{y} | \theta, M_l][\theta | M_l]P(M_l)}{\sum_{i=1}^4 [\mathbf{y} | \theta, M_i][\theta | M_i]P(M_i)} \quad (59)$$

for each model  $l = 1, \dots, 4$ .

- 7) Sample  $M_l^{(k+1)}$  from a categorical distribution with probabilities  $P(M_1 | \cdot)$ ,  $P(M_2 | \cdot)$ ,  $P(M_3 | \cdot)$ , and  $P(M_4 | \cdot)$ .
- 8) Increment the MCMC iteration  $t = t + 1$  and go to step 3.

A few of the terms in step 6 of the Barker and Link (2013) algorithm need further clarification with respect to the specific model set under consideration. The likelihood term for our Willow Tit occupancy model simplifies to  $[\mathbf{y} | \theta, M_l] \equiv [\mathbf{y} | p_t^{(t)}, \beta_{0,t}^{(t)}, \beta_t^{(t)}]$ , which can be found by integrating  $\mathbf{z}$  and  $\mathbf{v}$  out of the hierarchical model such that

$$[\mathbf{y} | p_t, \beta_{0,t}, \beta_t] = \prod_{i=1}^n \left( \psi_i p^{y_i} (1-p)^{J_i - y_i} I_{\{y_i > 0\}} \right) + \left( 1 - \psi_i + \psi_i (1-p)^{J_i} \right) I_{\{y_i = 0\}} \quad (60)$$

where we have omitted the MCMC indexing for clarity. In the integrated likelihood (Eq. 60),  $\psi_i = \phi(\mathbf{x}'_i \beta_t)$  and  $I_{\{\cdot\}}$  is an indicator variable that is one when the condition in the subscript is true and zero otherwise. The prior term can be factored into terms relevant for the current model being considered and terms for the remaining parameters:  $[\theta | M_l] \equiv [p_t^{(t)}][\beta_{0,t}^{(t)}][\beta_t^{(t)}][\mathbf{u}^{(t)}]$ . The last term,  $[\mathbf{u}^{(t)}]$ , is simply a product of independent standard normal distributions in our occupancy model.

This secondary MCMC algorithm required only seconds to run, as compared with the original model fits, which required minutes. Furthermore, we found the secondary MCMC algorithm suggested by Barker and Link (2013) easier to program than the inline RJMCMC algorithm because we didn't have to modify the actual model fitting code. Obtaining the posterior model probabilities from the secondary MCMC algorithm output simply requires calculating the number of times each model  $M_l^{(t)}$  is sampled out of the total number of MCMC iterations (e.g.,  $P(M_2 | \mathbf{y}) = 83\,200/160\,000 = 0.52$ ).

Several other alternatives exist for implementing RJMCMC and obtaining required BMA quantities. Notable among them are techniques for regression models that exploit orthogonality properties in the design matrix allowing for a simplification in the model sampler (Clyde et al. 1996). More recently, a form of data augmentation has been proposed to generalize these methods for cases where the design matrix is non-orthogonal (Ghosh and Clyde 2011). Overall, the suite of new approaches for model-based model selection is rapidly expanding and is making Bayesian model averaging more accessible than ever for ecologists. Still, fully automated software for performing BMA for a huge class of potential models is lacking due to the complexity of rigorously calculating the required quantities. As with many of the cutting-edge statistical methods, ecologists who wish to use them are acquiring the necessary statistical and computational skills to implement them on their own.

#### GUIDANCE

Thus far we have provided a fairly comprehensive review of methods for Bayesian model selection and multimodel inference, along with the advantages and

disadvantages of each. One can use this document as a reference in deciding what type of model selection is appropriate depending on the desired statistical inference in a particular project. Assuming that the researcher desires some form of model selection or multimodel inference, and that they plan to use Bayesian methods, we provide the following set of questions and answers to help guide the researcher in finding an appropriate set of tools:

- 1) Is the researcher planning a new study? If so, he or she may want to consider collecting two sets of data, one for training, and another for validation. When prediction is of utmost importance, there is no substitute for out-of-sample data in model selection. It may be time for a paradigm shift in the way we design ecological studies. If predictive model selection is desired, we need to collect data that facilitates inference on both parameters and models.
- 2) Is the researcher using a historical data set? (a) If the data set is large and computation time is not an overriding issue, the researcher may want to consider  $K$ -fold cross-validation for a set of candidate models or Bayesian regularization. Most Bayesian cross-validation implementations will require  $K$  separate fits of the model, thus increasing the computational time significantly. However, parallel computing is now possible on the desktop computer thanks to several user friendly software packages. So, cross-validation may not be as impractical as one might initially think. (b) If the data set is small,  $n$ -fold cross-validation over a set of candidate models or Bayesian regularization may be more appropriate. The caveat is that leave-one-out cross-validation is not as stable as  $K$ -fold for  $K < n$ . Small data sets are always going to present problems for statistical inference and there is not much one can do to alleviate these issues, regardless of statistical paradigm.
- 3) Is the researcher wanting to do prediction-based model selection with a simple Bayesian model when computational time is limited? If so, they might want to consider using DIC. As a prediction-based information criterion, DIC performs similar to AIC in choosing parsimonious models. The caveat is that, like AIC, DIC will also choose larger models than necessary when the sample size is large. The biggest caution about DIC arises when the posterior mean of the parameters does not describe the central tendency of the posterior distribution well. Thus, DIC is not appropriate when there exist multiple modes in the posterior. Furthermore, DIC is best as a selection criterion when the number of effective parameters is much smaller than the sample size, which may not be the case in hierarchical models where the number of latent variables scales with sample size.
- 4) Does the researcher want to do prediction-based model selection with a hierarchical Bayesian model when computational time is limited? If so, Gelman et al. (2014b) recommend using WAIC to select models. Unlike DIC, WAIC does not rely on posterior means of parameters, instead it uses the posterior predictive distribution and is the “most Bayesian” of all the information criteria. However, despite all the benefits of WAIC, it still only depends on within-sample data and its computationally friendly form requires an independence assumption at the data level, which is not appropriate for time series or spatial models. In these cases, posterior predictive loss provides an alternative.
- 5) Does the researcher desire model averaged inference on parameters or predictions? Bayes factors are the appropriate tool for doing Bayesian model averaging, but they often can only be approximated. Bayes factors can be approximated using BIC, but only under certain circumstances, and since BIC is not actually Bayesian, it has limited utility in a fully Bayesian setting. Hoeting et al. (1999) provided a good summary of methods for approximating model weights that have a formal justification. Note that, aside from BIC, none of the other information criteria have a solid foundation for Bayesian model averaging (e.g., AIC, DIC, WAIC). Bayes factors are not recommended in cases where models include improper priors (Spiegelhalter and Smith 1982).
- 6) Does the researcher want a fully integrated model fitting and selection procedure? If so, a model-based approach like indicator or Gibbs variable selection, stochastic search variable selection, or RJMCMC may be warranted. Furthermore, connections exist between many model-based approaches and BMA under certain conditions. These model-based methods perform best with some tuning of the algorithms, but when tuned, they perform quite well and seem to be more computationally efficient than cross-validation. As with information criteria, model-based model selection methods depend only on within-sample data and thus have the same set of caveats. Also, RJMCMC can be quite difficult to implement for certain models, but there are newer approaches that can be used to provide the same inference based on individual model fits (e.g., Barker and Link 2013).

#### CONCLUSION

Ecologists are fascinated with model selection, and many have customized their research questions around likelihood methods for model selection and multimodel inference as illustrated by the recent forum on  $P$  values and model selection in *Ecology* (2014, volume 95). Bayesian methods are becoming more common in ecological studies, but due to a fracturing of the literature pertaining to Bayesian model selection, it appears that many studies simply rely on conventional methods without much thought. Many Bayesian ecologists are aware of issues with certain Bayesian model

selection approaches (e.g., Bolker 2009), but are unaware of alternatives and how these alternatives may relate to each other. We have compiled and summarized the large body of literature on Bayesian model selection and multimodel inference methods in this guide so that ecologists can be better informed about their options.

What stands out to us is that, despite the seeming consensus among ecologists and wildlife biologists in how to perform model selection and multimodel inference, it is far from settled among statisticians; particularly in the Bayesian realm of inference. What also stands out is that nearly all model selection and multimodel inference methods are focused on improving predictive capabilities of models by balancing model fit and model parsimony. Prediction is often most important to the machine learning community (e.g., classification and regression trees, boosting and bagging algorithms) and related methods rely almost exclusively on out-of-sample data for model validation to improve prediction, but in the ecological and biological sciences, our scope seems to be limited to within-sample data. With an increasing ability to collect more data through, for example, better telemetry devices, remote sensing, citizen science efforts, and operations like NEON (National Ecological Observatory Network), ecologists are finally finding themselves with more data to answer scientific questions. Thus, model selection methods that rely on a separate set of validation data are now more accessible than ever for ecologists.

Cross-validation is an incredibly useful tool for model selection when only a single data set is available, a tool that is often overlooked or ignored on the grounds that it may be computationally infeasible. However, the current era of computing is seeing the most improvement in processor quantity and no longer in processor speed (Sutter 2005). The one thing that computers are getting better at is parallel processing, and that happens to strongly favor the notion of model selection via cross-validation. A bit of extra effort spent on bookkeeping aspects of programming can make true prediction-based model selection feasible through the parallelization of a cross-validation procedure. Using the occupancy model as an example, we demonstrated that parallel programming requires relatively little extra effort to implement but can improve computational efficiency dramatically (e.g., from hours to minutes, sometimes seconds).

When it seems that fitting a single model is the computational bottleneck, we need to remember that there are several entire subfields within statistics and computer science devoted to finding more efficient ways to specify and fit models. Automated MCMC software has been a boon for science, allowing ecologists to easily specify and fit complicated Bayesian models (e.g., Kery 2010), but a common complaint is that these software packages are slow. Fortunately, a wave of new automatic Bayesian software is becoming available (e.g., INLA, STAN, LibBi) that has shown dramatic

increases in speed, but improvements can also be gained just by creating our own MCMC algorithms. This gives us the flexibility to use model reparameterizations and newer computational tricks such as variational Bayes (e.g., Omerod and Wand 2010) and statistical emulators (e.g., Hooten et al. 2011) to speed up the model fitting process, which in turn aids in out-of-sample model selection.

Finally, as a closing thought, we feel that it is the right time for ecologists to become more open-minded about the use of strong priors. It is somewhat ironic that many popular non-Bayesian statistical methods (e.g., model selection, penalized likelihood, Lasso) depend on the implicit use of strong priors while at the same time Bayesians are warned against them. Bayesian priors provide a formal mechanism for placing constraints on models and, when used correctly, such constraints can be incredibly helpful (e.g., Moreno and Lele 2010). Furthermore, seemingly vague priors can have a dubious effect on inference (Seaman et al. 2012) in models commonly used in ecological analyses. Yet, stronger priors can help with model selection, multicollinearity, and algorithm stability, not to mention formally incorporating existing scientific information into new analyses (e.g., Garrard et al. 2012).

#### ACKNOWLEDGMENTS

The authors thank two anonymous reviewers as well as several colleagues, including David Anderson, Brian Brost, Frances Buderman, Ken Burnham, Alison Cartwright, Bob Dorazio, Brian Gerber, Tabitha Graves, Ephraim Hanks, Megan Higgs, Jennifer Hoeting, Devin Johnson, Shannon Kay, Bill Link, Kiona Ogle, Ann Raiho, Jay Rotella, Andy Royle, Viviana Ruiz-Gutierrez, Maria Uriarte, Jay Ver Hoef, and Perry Williams for valuable insight and early discussions about this work. Support for this work was provided by NSF 000347455 and NSF EF1241856. Any use of trade, firm, or product names is for descriptive purposes only and does not imply endorsement by the U.S. Government.

#### LITERATURE CITED

- Albert, J., and S. Chib. 1990. Bayesian analysis of binary and polychotomous response data. *Journal of the American Statistical Association* 88:669–679.
- Ando, T., and R. Tsay. 2010. Predictive likelihood for Bayesian model selection and averaging. *International Journal of Forecasting* 26:744–763.
- Barker, R. J., and W. A. Link. 2013. Bayesian multimodel inference by RJMCMC: A Gibbs sampling approach. *American Statistician* 67:150–156.
- Berger, J. O. 2006. *Statistical decision theory and Bayesian analysis*. Springer, New York, New York, USA.
- Bernardo, J. M. 1979. Expected information as expected utility. *Annals of Statistics* 7:686–690.
- Bernardo, J. M., and A. F. M. Smith. 1994. *Bayesian theory*. John Wiley, New York, New York, USA.
- Bolker, B. 2008. *Ecological models and data* in R. Princeton University Press, Princeton, New Jersey, USA.
- Bolker, B. 2009. Learning hierarchical models: advice for the rest of us. *Ecological Applications* 19:588–592.
- Bondell, H. D., and B. J. Reich. 2012. Consistent high-dimensional Bayesian variable selection via penalized credible regions. *Journal of the American Statistical Association* 107:1610–1624.



- Burnham, K. P., and D. R. Anderson. 2002. Model selection and multimodel inference. Second edition. Springer-Verlag, Berlin, Germany.
- Carlin, B. R., and S. Chib. 1995. Bayesian model choice via Markov Chain Monte Carlo methods. *Journal of the Royal Statistical Society B* 57:473–484.
- Celeux, G., F. Forbes, C. P. Robert, and D. M. Titterton. 2006. Deviance information criteria for missing data models. *Bayesian Analysis* 1:651–674.
- Clark, J. S. 2005. Why environmental scientists are becoming Bayesians. *Ecology Letters* 8:2–14.
- Clark, J. S. 2007. Models for ecological data: an introduction. Princeton University Press, Princeton, New Jersey, USA.
- Clyde, M. A., H. Desimone, and G. Parmigiani. 1996. Prediction via orthogonalized model mixing. *Journal of the American Statistical Association* 91:1197–1208.
- Congdon, P. 2006. Bayesian model choice based on Monte Carlo estimates of posterior model probabilities. *Computational Statistics and Data Analysis* 50:346–357.
- Cressie, N., C. A. Calder, J. S. Clark, J. M. Ver Hoef, and C. K. Wikle. 2009. Accounting for uncertainty in ecological analysis: the strengths and limitations of hierarchical statistical modeling. *Ecological Applications* 19:553–570.
- Czado, C., T. Gneiting, and L. Held. 2009. Predictive model assessment for count data. *Biometrics* 65:1254–1261.
- Dahlgren, J. P. 2010. Alternative regression methods are not considered in Murtaugh (2009) or by ecologists in general. *Ecology Letters* 13:E7–E9.
- Dellaportas, P., J. J. Forster, and I. Ntzoufras. 1997. On Bayesian model and variable selection using MCMC. Technical report. Department of Statistics, Athens University of Economics and Business, Athens, Greece.
- Dorazio, R. M., M. Kery, J. A. Royle, and M. Plattner. 2010. Models for inference in dynamic metacommunity systems. *Ecology* 91:2466–2475.
- Dorazio, R. M., and D. T. Rodriguez. 2012. A Gibbs sampler for Bayesian analysis of site-occupancy data. *Methods in Ecology and Evolution* 3:1093–1098.
- Garrard, G. E., M. A. McCarthy, P. A. Vesik, J. Q. Radford, and A. F. Bennett. 2012. A predictive model of avian natal dispersal distance provides prior information for investigating response to landscape change. *Journal of Animal Ecology* 81:14–23.
- Geisser, S. 1993. Predictive inference: an introduction. Chapman and Hall, London, UK.
- Gelfand, A. E., and S. K. Ghosh. 1998. Model choice: a minimum posterior predictive loss approach. *Biometrika* 85:1–13.
- Gelfand, A. E., and A. F. M. Smith. 1990. Sampling-based approaches to calculating marginal densities. *Journal of the American Statistical Association* 85:398–409.
- Gelman, A., J. B. Carlin, H. S. Stern, D. B. Dunson, A. Vehtari, and D. B. Rubin. 2014a. Bayesian data analysis. Third edition. Chapman and Hall/CRC, Boca Raton, Florida, USA.
- Gelman, A., J. Huang, and A. Vehtari. 2014b. Understanding predictive information criteria for Bayesian models. *Statistics and Computing*. <http://dx.doi.org/10.1007/s11222-013-9416-2>
- Gelman, A., and C. R. Shalizi. 2012. Philosophy and the practice of Bayesian statistics. *British Journal of Mathematical and Statistical Psychology* 66:8–38.
- George, E. I., and R. E. McCulloch. 1993. Variable selection via Gibbs sampling. *Journal of the American Statistical Association* 85:398–409.
- Ghosh, J., and M. A. Clyde. 2011. Rao-Blackwellization for Bayesian variable selection and model averaging in linear and binary regression: a novel data augmentation approach. *Journal of the American Statistical Association* 106:1041–1052.
- Gneiting, T. 2011. Making and evaluating point forecasts. *Journal of the American Statistical Association* 106:746–762.
- Gneiting, T., and A. E. Raftery. 2007. Strictly proper scoring rules, prediction, and estimation. *Journal of the American Statistical Association* 102:359–378.
- Godsill, S. J. 2001. On the relationship between Markov chain Monte Carlo methods for model uncertainty. *Journal of Computational and Statistical Graphics* 10:230–248.
- Gotelli, N. J., and A. M. Ellison. 2012. A primer of ecological statistics. Second edition. Sinauer Associates, Sunderland, Massachusetts, USA.
- Graham, M. H. 2003. Confronting multicollinearity in ecological multiple regression. *Ecology* 84:2809–2815.
- Green, P. J. 1995. Reversible jump Markov chain Monte Carlo computation and Bayesian model determination. *Biometrika* 82:711–732.
- Hastie, D. I., and P. J. Green. 2012. Model choice using reversible jump Markov chain Monte Carlo. *Statistica Neerlandica* 66:309–338.
- Hastie, T., R. Tibshirani, and J. Friedman. 2009. Elements of statistical learning: data mining, inference, and prediction. Second edition. Springer, New York, New York, USA.
- Held, L., B. Schrodle, and H. Rue. 2010. Posterior and cross-validatory predictive checks: a comparison of MCMC and INLA. Pages 91–109 in T. Kneib and G. Tutz, editors. *Statistical modelling and regression structures: festschrift in honour of Ludwig Fahrmeir*. Springer, New York, New York, USA.
- Hobbs, N. T. 2009. New tools for insight from ecological models and data. *Ecological Applications* 19:551–552.
- Hobbs, N. T., and M. B. Hooten. *In press*. Bayesian models: a statistical primer for ecologists. Princeton University Press, Princeton, New Jersey, USA.
- Hoeting, J. A., D. Madigan, A. E. Raftery, and C. T. Volinsky. 1999. Bayesian model averaging: a tutorial. *Statistical Science* 14:382–417.
- Hooten, M. B., D. R. Larsen, and C. K. Wikle. 2003. Predicting the spatial distribution of ground flora on large domains using a hierarchical Bayesian model. *Landscape Ecology* 18:487–502.
- Hooten, M. B., W. B. Leeds, J. Fiechter, and C. K. Wikle. 2011. Assessing first-order emulator inference for physical parameters in nonlinear mechanistic models. *Journal of Agricultural, Biological and Environmental Statistics* 16:475–494.
- Jeffreys, H. 1961. Theory of probability. Third edition. Oxford University Press, Oxford, UK.
- Johnson, D. S., P. B. Conn, M. B. Hooten, J. Ray, and B. Pond. 2013. Spatial occupancy models for large data sets. *Ecology* 94:801–808.
- Johnson, D. S., and J. A. Hoeting. 2011. Bayesian multimodel inference for geostatistical regression models. *PLoS ONE* 6:e25677.
- Johnson, J. B., and K. S. Omland. 2004. Model selection in ecology and evolution. *Trends in Ecology and Evolution* 19:101–108.
- Kass, R. E., and A. E. Raftery. 1995. Bayes factors. *Journal of the American Statistical Association* 90:773–795.
- Kery, M. 2010. Introduction to WinBUGS for ecologists. Academic Press, Burlington, Massachusetts, USA.
- Kery, M., and H. Schmid. 2004. Monitoring programs need to take into account imperfect species detectability. *Basic and Applied Ecology* 5:65–73.
- Knaus, J. 2013. snowfall: easier cluster computing (based on snow). R package version 1.84-4. <http://cran.r-project.org/package=snowfall>
- Kuo, L., and B. Mallick. 1998. Variable selection for regression models. *Sankhya, Series B* 60:65–81.

- Kyung, M., J. Gill, M. Ghosh, and G. Casella. 2010. Penalized regression, standard errors, and Bayesian lassos. *Bayesian Analysis* 5:369–412.
- Laud, P., and J. Ibrahim. 1995. Predictive model selection. *Journal of the Royal Statistical Society B* 57:247–262.
- Lehmann, E. L., and G. Casella. 1998. *Theory of point estimation*. Springer, New York, New York, USA.
- Link, W. A., and R. J. Barker. 2006. Model weights and the foundations of multimodel inference. *Ecology* 87:2626–2635.
- Link, W. A., and R. J. Barker. 2010. *Bayesian inference: with ecological applications*. Academic Press, London, UK.
- MacKenzie, D. I., J. D. Nichols, J. E. Hines, M. G. Knutson, and A. B. Franklin. 2003. Estimating site occupancy, colonization, and local extinction when a species is detected imperfectly. *Ecology* 84:2200–2255.
- MacKenzie, D. I., J. D. Nichols, G. B. Lachman, S. Droege, J. A. Royle, and C. A. Langtimm. 2002. Estimating site occupancy rates when detection probabilities are less than one. *Ecology* 83:2248–2255.
- MacKenzie, D. I., J. D. Nichols, J. A. Royle, K. H. Pollock, L. L. Bailey, and J. E. Hines. 2006. *Occupancy estimation and modeling*. Elsevier, Amsterdam, The Netherlands.
- Madigan, D., and A. E. Raftery. 1994. Model selection and accounting for model uncertainty in graphical models using Occam's window. *Journal of the American Statistical Association* 89:1535–1546.
- Martin, T. G., B. A. Wintle, J. R. Rhodes, P. M. Kuhnert, S. A. Field, S. J. Low-Choy, A. J. Tyre, and H. Possingham. 2005. Zero-tolerance in ecology: improving ecological inference by modelling the source of zero observations. *Ecology Letters* 8:1235–1246.
- Miller, A. 2002. *Subset selection in regression*. Chapman and Hall/CRC, Boca Raton, Florida, USA.
- Moreno, M., and S. R. Lele. 2010. Improved estimation of site occupancy using penalized likelihood. *Ecology* 91:341–346.
- O'Hara, R. B., and M. J. Sillanpaa. 2009. A review of Bayesian variable selection methods: What, how and which. *Bayesian Analysis* 4:85–118.
- Omerod, J. T., and M. P. Wand. 2010. Explaining variational approximations. *American Statistician* 64:140–153.
- Park, T., and G. Casella. 2008. The Bayesian Lasso. *Journal of the American Statistical Association* 103:681–686.
- Pettit, L. I. 1990. The conditional predictive ordinate for the normal distribution. *Journal of the American Statistical Association* 52:175–184.
- Plummer, M. 2002. Discussion of the paper by Spiegelhalter et al. *Journal of the Royal Statistical Society B* 64:620.
- Plummer, M. 2008. Penalized loss functions for Bayesian model comparison. *Biostatistics* 9:523–539.
- R Core Team. 2013. *R: a language and environment for statistical computing*. R Foundation for Statistical Computing, Vienna, Austria. <http://www.r-project.org>
- Richardson, S. 2002. Discussion of the paper by Spiegelhalter et al. *Journal of the Royal Statistical Society B* 64:626–227.
- Ripley, B. D. 1996. *Pattern recognition and neural networks*. Cambridge University Press, Cambridge, UK.
- Royle, J. A., and R. M. Dorazio. 2008. *Hierarchical modeling and inference in ecology*. Academic Press.
- Schwarz, G. E. 1978. Estimating the dimension of a model. *Annals of Statistics* 6:461–464.
- Seaman, J. W., III, J. W. Seaman, Jr., and J. D. Stamey. 2012. Hidden dangers of specifying noninformative priors. *American Statistician* 66:77–84.
- Spiegelhalter, D. J., N. G. Best, B. P. Carlin, and A. van der Linde. 2002. Bayesian measures of model complexity and fit. *Journal of the Royal Statistical Society B* 64:583–639.
- Spiegelhalter, D. J., and A. F. M. Smith. 1982. Bayes factors for linear and log-linear models with vague prior information. *Journal of the Royal Statistical Society B* 44:377–387.
- Stone, M. 1977. An asymptotic equivalence of choice of model cross-validation and Akaike's criterion. *Journal of the Royal Statistical Society B* 36:44–47.
- Sutter, H. 2005. The free lunch is over: a fundamental turn toward concurrency in software. *Dr. Dobbs Report* 30(3).
- Tanner, M. A. 1996. *Tools for statistical inference: methods for the exploration of posterior distributions and likelihood functions*. Third edition. Springer, New York, New York, USA.
- Tibshirani, R. 1996. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society B* 58:267–288.
- Vehtari, A., and J. Ojanen. 2012. A survey of Bayesian predictive methods for model assessment, selection and comparison. *Statistics Surveys* 6:142–228.
- Watanabe, S. 2010. Asymptotic equivalence of Bayes cross-validation and widely applicable information criterion in singular learning theory. *Journal of Machine Learning Research* 11:3571–3594.
- Watanabe, S. 2013. A widely applicable Bayesian information criterion. *Journal of Machine Learning Research* 14:867–897.

## SUPPLEMENTAL MATERIAL

## Ecological Archives

The Supplement is available online: <http://dx.doi.org/10.1890/14-0661.1.sm>