

Chapter 1

Introduction to Bayesian Analysis of GL(M)Ms Using R/WinBUGS

A major theme of this book is that spatial capture-recapture models are, for the most part, just generalized linear models (GLMs) wherein the covariate, distance between trap and home range center, is partially or fully unobserved – and therefore regarded as a random effect. Such models are usually referred to as Generalized Linear Mixed Models (GLMMs) and, therefore, SCR models can be thought of as a specialized type of GLMM. Naturally then, we should consider analysis of these slightly simpler models in order to gain some experience and, hopefully, develop a better understanding of spatial capture-recapture models.

In this chapter, we consider classes of GLM models - Poisson and binomial (i.e., logistic regression) GLMs - that will prove to be enormously useful in the analysis of capture-recapture models of all kinds. Many readers are probably familiar with these models because they represent probably the most generally useful models in all of Ecology and, as such, have received considerable attention in many introductory and advanced texts. We focus on them here in order to introduce the readers to the analysis of such models in **R** and **WinBUGS**, which we will translate directly to the analysis of SCR models in subsequent chapters.

Bayesian analysis is convenient for analyzing GLMMs because it allows us to work directly with the conditional model – i.e., the model that is conditional on the random effects, using computational methods known as Markov chain Monte Carlo (MCMC). Learning how to do Bayesian analysis of GLMs and GLMMs in **WinBUGS** is, in part, the purpose of this chapter. While we use **WinBUGS** to do the Bayesian computations, we organize and summarize our data and execute **WinBUGS** from within **R** using the useful package **R2WinBUGS**

(Sturtz et al., 2005). Kéry (2010), and Kéry and Schaub (2011) provide excellent introductions to the basics of Bayesian analysis and GLMs at an accessible level. We don't want to be too redundant with those books and so we avoid a detailed treatment of Bayesian methodology - instead just providing a cursory overview so that we can move on and attack the problems we're most interested in related to spatial capture-recapture. In addition, there are a number of texts that provide general introductions to Bayesian analysis, MCMC, and their applications in Ecology including McCarthy (2007), Kéry (2010), Link and Barker (2009), and King (2009).

While this chapter is about Bayesian analysis of GLMMs, such models are routinely analyzed using likelihood methods too, as discussed by Royle and Dorazio (2008), and Kéry (2010). Indeed, likelihood analysis of such models is the primary focus of many applied statistics texts, a good one being Zuur et al. (2009). Later in this book, we will use likelihood methods to analyze SCR models but, for now, we concentrate on providing a basic introduction to Bayesian analysis because that is the approach we will use in a majority of cases in later chapters.

1.1 Notation

We will sometimes use conventional “bracket notation” to refer to probability distributions. If y is a random variable the $[y]$ indicates its distribution or its probability density/mass function (pdf, pmf) depending on context. If x is another random variable then $[y|x]$ is the conditional distribution of y given x , and $[y, x]$ is the joint distribution of y and x . To differentiate specific distributions in some contexts we might label them $g(y)$, $g(y|\theta)$, $f(x)$, or similar. We will also write $y \sim \text{Normal}(\mu, \sigma^2)$ to indicate that y “is distributed as” a normal random variable with parameters μ and σ^2 . The expected value or mean of a random variable is $E[y] = \mu$, and $\text{Var}[y] = \sigma^2$ is the variance of y . To indicate specific observations we'll use an index such as “ i ”. So, y_i for $i = 1, 2, \dots, n$ indicates observations for n individuals. Finally, we write $\text{Pr}(y)$ to indicate specific probabilities, i.e., of events “ y ” or similar.

To illustrate these concepts and notation, suppose z is a binary outcome (e.g., species occurrence) and we might assume the model: $z \sim \text{Bern}(p)$ for observations. Under this model $\text{Pr}(z = 1) = \psi$, which is also the expected value $E[z] = \psi$. The variance is $\text{Var}[z] = \psi * (1 - \psi)$ and the probability mass function (pmf) is $[z] = \psi^z (1 - \psi)^{1-z}$. Sometimes we write $[z|\psi]$ when it is important to emphasize the conditional dependence of z on ψ . As another example, suppose y is a random variable denoting whether or not a species is detected if an occupied site is surveyed. In this case it might be natural to express the pmf of the observations y conditional on z . That is, $[y|z]$. In this case, $[y|z = 1]$ is the conditional pmf of y given that a site is occupied, and it is natural to assume that $[y|z = 1] = \text{Bern}(p)$ where p is the “detection probability” - the probability that we detect the species, given that it is present. The model for the observations y is completely specified once we describe the other conditional

pmf $[y|z = 0]$. For this conditional distribution it is sometimes reasonable to assume $\Pr(y = 1|z = 0) = 0$ (MacKenzie et al. (2002); see also Royle and Link (2006)). That is, if the species is absent, the probability of detection is 0. This implies that $\Pr(y = 0|z = 0) = 1$. To allow for situations in which the true state z is unobserved, we assume that $[z]$ is Bernoulli with parameter ψ . In this case, the marginal distribution of y is

$$[y] = [y|z = 1]Pr(z = 1) + [y|z = 0]Pr(z = 0)$$

because $[y|z = 0]$ is a point mass at $y = 0$, by assumption, then

$$\Pr(y = 1) = p\psi$$

And

$$\Pr(y = 0) = (1 - p) * \psi + (1 - \psi)$$

1.2 GLMs and GLMMs

We have asserted already that SCR models work out most of the time to be variations of GLMs and GLMMs. Some of you might therefore ask: What are GLMs and GLMMs, anyhow? These models are covered extensively in many very good applied statistics books and we refer the reader elsewhere for a detailed introduction. We think Kéry (2010), Kéry and Schaub (2011), and Zuur et al. (2009) are all accessible treatments of considerable merit. Here, we'll give the 1 minute treatment of GLMMs, not trying to be complete but rather only to preserve a coherent organization to the book.

The generalized linear model (GLM) is an extension of standard linear models by allowing the response variable to have some distribution from the exponential family of distributions (i.e., not just normal). This includes the normal distribution but also dozens of others such as the Poisson, binomial, gamma, exponential, and many more. In addition, GLMs allow the response variable to be related to the predictor variables (i.e., covariates) using a link function, which is usually nonlinear. Finally, GLMs typically accommodate a relationship between the mean and variance. The classical reference for GLMs is Nelder and Wedderburn (1972) and also McCullagh and Nelder (1989). The GLM consists of three components:

1. A probability distribution for the dependent variable y , from a class of probability distributions known as the exponential family.
2. A "linear predictor" $\eta = \mathbf{X}\beta$.
3. A link function g that relates $E[y]$ to the linear predictor, $E[y] = \mu = g^{-1}(\eta)$. Therefore $g(E[y]) = \eta$.

The dependent variable y is assumed to be an outcome from a distribution of the exponential family which includes many common distributions including

the normal, gamma, Poisson, binomial, and many others. The mean of the distribution of y is assumed to depend on predictor variables x according to

$$g(E[y]) = \mathbf{x}'\beta$$

where $E[y]$ is the expected value of y , and $\mathbf{x}'\beta$ is termed the *linear predictor*, i.e., a linear function of the predictor variables with unknown parameters β to be estimated. The function g is the link function. In standard GLMs, the variance of y is a function V of the mean of y : $Var(y) = V(\mu)$ (see below for examples).

A Poisson GLM posits that $y \sim \text{Poisson}(\lambda)$ with $E[y] = \lambda$ and usually the model for the mean is specified using the *log link function* by

$$\log(\lambda_i) = \beta_0 + \beta_1 * x_i$$

The variance function is $V(y_i) = \lambda_i$. The binomial GLM posits that $y_i \sim \text{Binomial}(K, p)$ where K is the fixed sample size parameter and $E[y_i] = K * p_i$. Usually the model for the mean is specified using the *logit link function* according to

$$\text{logit}(p_i) = \beta_0 + \beta_1 * x_i$$

Where $\text{logit}(u) = \log(u/(1-u))$. The inverse-logit function, g^{-1} , is a function we will refer to as “expit”, so that $\text{expit}(u) = \exp(u)/(1 + \exp(u))$.

A GLMM is the extension of GLMs to accommodate “random effects”. Often this involves adding a normal random effect to the linear predictor, and so a simple example is:

$$\log(\lambda_i) = \alpha_i + \beta_1 * x_i$$

where

$$\alpha_i \sim \text{Normal}(\mu, \sigma^2)$$

1.3 Bayesian Analysis

Bayesian analysis is unfamiliar to many ecological researchers because older cohorts of ecologists were largely educated in the classical statistical paradigm of frequentist inference. But advances in technology and increasing exposure to benefits of Bayesian analysis are fast making Bayesians out of people or at least making Bayesian analysis an acceptable, general, alternative to classical, frequentist inference.

Conceptually, the main thing about Bayesian inference is that it uses probability directly to characterize uncertainty about things we don’t know. “Things”, in this case, are parameters of models and, just as it is natural to characterize uncertain outcomes of stochastic processes using probability, it seems natural also to characterize information about unknown “parameters” using probability. At least this seems natural to us and, we think, most ecologists either explicitly adopt that view or tend to fall into that point of view naturally. Conversely, frequentists use probability in many different ways, but never to characterize

uncertainty about parameters¹ Instead, frequentists use probability to characterize the behavior of *procedures* such as estimators or confidence intervals (see below), which can lead to some inelegant or unnatural interpretations of things. It is paradoxical that people readily adopt a philosophy of statistical inference in which the things you don't know (i.e., parameters) should *not* be regarded as random variables, so that, as a consequence, one cannot use probability to characterize one's state of knowledge about them.

1.3.1 Bayes Rule

As its name suggests, Bayesian analysis makes use of Bayes' rule in order to make direct probability statements about model parameters. Given two random variables z and y , Bayes rule relates the two conditional probability distributions $[z|y]$ and $[y|z]$ by the relationship:

$$[z|y] = [y|z][z]/[y]$$

Bayes' rule itself is a mathematical fact and there is no debate in the statistical community as to its validity and relevance to many problems. Generally speaking, these distributions are characterized as follows: $[y|z]$ is the conditional probability distribution of y *given* z , $[z]$ is the marginal distribution of z and $[y]$ is the marginal distribution of y . In the context of Bayesian inference we usually associate specific meanings in which $[y|z]$ is thought of as "the likelihood", $[z]$ as the "prior" and so on. We leave this for later because here the focus is on this expression of Bayes rule as a basic fact of probability.

As an example of a simple application of Bayes rule, consider the problem of determining species presence at a sample location based on imperfect survey information. Let z be a binary random variable that denotes species presence ($z = 1$) or absence ($z = 0$), let $\Pr(z = 1) = \psi$ where ψ is usually called occurrence probability, "occupancy" (MacKenzie et al., 2002) or "prevalence". Let y be the *observed* presence ($y = 1$) or absence ($y = 0$), and let p be the probability that a species is detected in a single survey at a site given that it is present. Thus, $\Pr(y = 1|z = 1) = p$. The interpretation of this is that, if the species is present, we will only observe presence with probability p . In addition, we assume here that $\Pr(y = 1|z = 0) = 0$. That is, the species cannot be detected if it is not present which is a conventional view adopted in most biological sampling problems (but see Royle and Link (2006)). If we survey a site T times but never detect the species, then this clearly does not imply that the species is not present ($z = 0$) at this site. Rather, our degree of belief in $z = 0$ should be made with a probabilistic statement $\Pr(z = 1|y_1 = 0, \dots, y_T = 0)$. If the T surveys are independent so that we might regard y_t as *iid* Bernoulli trials, then the total number of detections, say y , is Binomial with probability p then we can use Bayes rule to compute the probability that it is present given that

¹To hear this will be shocking to some readers perhaps.

177 it is not detected in T samples. In words, the expression we seek is:

$$\Pr(\text{present}|\text{not detected}) = \frac{\Pr(\text{not detected}|\text{present}) \Pr(\text{present})}{\Pr(\text{detected})}$$

178 Mathematically, this is

$$\Pr(z = 1|y = 0) = \Pr(y = 0|z = 1) \Pr(z = 1) / \Pr(y = 0) = [(1-p)^T \psi] / [(1-p)^T \psi + (1-\psi)].$$

179 To apply this, suppose that $T = 2$ surveys are done at a wetland for a species
 180 of frog, and the species is not detected there. Suppose further that $\psi = .8$ and
 181 $p = .5$ are obtained from a prior study. Then the probability that the species is
 182 present at this site is $.25 * .8 / (.25 * .8 + .2) = 0.50$. That is, there seems to be
 183 about a 50/50 chance that the site is occupied despite the fact that the species
 184 wasn't observed there.

185 In summary, Bayes' rule provides a simple linkage between the conditional
 186 probabilities $[y|z]$ and $[z|y]$ which is useful whenever one needs to deduce one
 187 from the other. Bayes' rule as a basic fact of probability is not disputed.

188 1.3.2 Bayesian Inference

189 What is controversial to some is the scope and manner in which Bayes rule is
 190 applied by Bayesian analysts. Bayesian analysts assert that Bayes rule is rele-
 191 vant, in general, to all statistical problems by regarding all unknown quantities
 192 of a model as realizations of random variables - this includes "data", latent
 193 variables, and also "parameters". Classical (non-Bayesian) analysts sometimes
 194 object to regarding "parameters" as outcomes of random variables. Classically,
 195 parameters are thought of as "fixed but unknown" (using the terminology of
 196 classical statistics). Of course, in Bayesian analysis they are also unknown
 197 and, in fact, there is a single data-generating value and so they are also fixed.
 198 The difference is that this fixed but unknown value is regarded as having been
 199 generated from some probability distribution. Specification of that probability
 200 distribution is necessary to carryout Bayesian analysis, but it is not required in
 201 classical frequentist inference.

202 To see the general relevance of Bayes rule in the context of statistical infer-
 203 ence, let y denote observations - i.e., "data" - and let $[y|\theta]$ be the observation
 204 model (often colloquially referred to as the "likelihood"). Suppose θ is a
 205 parameter of interest having (prior) probability distribution $[\theta]$. These are com-
 206 bined to obtain the posterior distribution using Bayes' rule, which is:

$$[\theta|y] = [y|\theta][\theta]/[y]$$

207 Asserting the general relevance of Bayes rule to all statistical problems, we
 208 can conclude that the two main features of Bayesian inference are that: (1)
 209 "parameters" θ are regarded as realizations of a random variable and, as a
 210 result, (2) inference is based on the probability distribution of the parameters
 211 given the data, $[\theta|y]$, which is called the posterior distribution. This is the

result of using Bayes rule to combine “the likelihood” and the prior distribution. The key concept is regarding parameters as realizations of a random variable because, once you admit this conceptual view, this leads directly to the posterior distribution, a very natural quantity upon which to base inference about things we don’t know - including parameters of statistical models. In particular, $[\theta|y]$ is a probability distribution for θ and therefore we can make direct probability statements to characterize uncertainty about θ .

The denominator of our invocation of Bayes rule, $[y]$, is the marginal distribution of the data y . We note without further remark right now that, in many practical problems, this can be an enormous pain to compute. The main reason that the Bayesian paradigm has become so popular in the last 20 years or so is because methods exist for characterizing the posterior distribution that do not require that we possess a mathematical understanding of $[y]$, i.e., we never have to compute it or know what it looks like, or know anything specific about it.

A common misunderstanding on the distinction between Bayesian and frequentist inference goes something like this “in frequentist inference parameters are fixed but unknown but in a Bayesian analysis parameters are random.” At best this is a sad caricature of the distinction and at worst it is downright wrong. What is true is that, to a Bayesian, parameters are random variables. However, a Bayesian assumes, just like a frequentist, that there was a single data-generating value of that parameter - a fixed, and unknown value that produced the given data set. The distinction between Bayesian and frequentist approaches is that Bayesians regard the parameter as a random variable, and its value as the outcome of a random value, on par with the observations. This allows Bayesians to use probability to make direct probability statements about parameters. Frequentist inference procedures do not permit direct probability statements to be made about parameter values – because parameters are not random variables!

While we can understand the conceptual basis of Bayesian inference merely by understanding Bayes rule – that’s really all there is to it – it is not so easy to understand the basis of classical “frequentist” inference which is mostly like² a “basket of methods” with little coherent organization. What is mostly coherent in frequentist inference is the manner in which items in this basket of methods are evaluated – the performance of a given procedure is evaluated by “averaging over” hypothetical realizations of y , regarding the *estimator* as a random variable. For example, if $\hat{\theta}$ is an estimator of θ then the frequentist is interested in $E_y[\hat{\theta}|y]$ which is used to characterize bias. If the expected value of $\hat{\theta}$, when averaged over realizations of y , is equal to θ , then $\hat{\theta}$ is unbiased.

The view of parameters as fixed constants and estimators as random variables leads to interpretations that are not so straightforward. For example confidence intervals having the interpretation “95% probability that the interval contains the true value” and p-values being “the probability of observing an outcome as extreme or more than the one observed.” These are far from intuitive interpretations to most people. Moreover, this is conceptually prob-

²Characterization from Sims REF XYZ

blematic to some because the hypothetical realizations that characterize the performance of our procedure we will never get to observe.

While we do tend to favor Bayesian inference for the conceptual simplicity (parameters are random, posterior inference), we mostly advocate for a pragmatic non-partisanship approach to inference because, frankly, some of these “bucket of methods” are actually very convenient in certain situations as we will see in later chapters.

1.3.3 Prior distributions

The prior distribution $[\theta]$ is an important feature of Bayesian inference. As a conceptual matter, the prior distribution characterizes “prior beliefs” or “prior information” about a parameter. Indeed, an oft-touted benefit of Bayesian analysis is the ease with which prior information can be included in an analysis. However, more commonly, the prior is chosen to express a lack of prior information, even if previous studies have been done and even if the investigator does in fact know quite a bit about a parameter. This is because the manner in which prior information is embodied in a prior (and the amount of information) is usually very subjective and thus the result can wind up being very contentious, e.g., if different investigators might report different results based on subjective assessments of things. Thus it is usually better to “let the data speak” and use priors that reflect absence of information beyond the data set being analyzed.

But still the need occasionally arises to embody prior information or beliefs about a parameter formally into the estimation scheme. In SCR models we often have a parameter that is closely linked to “home range radius” and thus auxiliary information on the home range size of a species can be used as prior information (e.g., see Chandler and Royle (2012) ; also chapter XYZ).

XXXXXXXX

noninformative prior on one scale is informative on another scale. e.g., flat prior on $\logit(p)$ is very different from $\text{uniform}(0,1)$ on p ... show graphic.....

reference to non-invariance of prior distributions to transformation.....

XXXXXXXX

1.3.4 Posterior Inference

In Bayesian inference, we are not focusing on estimating a single point or interval but rather on characterizing a whole distribution – the posterior distribution – from which one can report any summary of interest. A point estimate might be the posterior mean, median, mode, etc.. In many applications in this book, we will compute 95% Bayesian intervals using the 2.5% and 97.5% quantiles of the posterior distribution. For such intervals, it is correct to say $\Pr(L < \theta < U) = 0.95$. That is, “the probability that θ is between L and U is 0.95”. It is not a subtle thing that this cannot be said using frequentist methods - although people tend to say it anyway and not really understand why it is wrong or even that it is wrong. This is actually a failing of frequentist ideas and the inability of frequentists to get people to overcome their natural tendency to use probability

298 - which is something that, as a frequentist, you simply cannot do in the manner
 299 that you would like to.

300 Posterior inference is the main practical element of Bayesian analysis. We
 301 get to make an inference conditional on the data that we actually observed -
 302 i.e., what we actually know. To us, this seems logical - to condition on what
 303 we know. Conversely, frequentist inference is based on considering average per-
 304 formance over hypothetical unobserved data sets (i.e., the “relative frequency”
 305 interpretation of probability). Frequentists know that their procedures work
 306 well when averaged over all hypothetical, unobserved, data sets but no one ever
 307 really knows how well they work for the specific data set analyzed. That seems
 308 like a relevant question to biologists who oftentimes only have their one, ex-
 309 tremely valuable, data set. This distinction comes into play a lot in exposing
 310 philosophical biases in the peer review of statistical analyses in ecology in the
 311 sense that, despite these opposing conceptual views to inference (i.e. condi-
 312 tional on the data you have, or averaged over hypothetical realizations), those
 313 who conduct a Bayesian analysis are often (in ecology, almost always) required
 314 to provide a frequentist evaluation of their Bayesian procedure.

315 1.3.5 Small sample inference

316 Using Bayesian inference, we obtain an estimate of the posterior distribution
 317 which is an exhaustive summary of the state-of-knowledge about an unknown
 318 quantity. It is the posterior distribution - not an estimate of that thing. It is
 319 also not, usually, an approximation except to within Monte Carlo error (in cases
 320 where we use simulation to calculate it). One of the great virtues of Bayesian
 321 analysis which is not really appreciated is that it is completely valid for any
 322 particular sample size. i.e., it is $[\theta|y]$, as precise as we claim it to be based on
 323 our ability to do calculations, for the particular sample size and observations
 324 that we have even if we have only a single datum y . The same cannot be said
 325 for almost all frequentist procedures in which estimates or variances are very
 326 often (almost always in practice) based on “asymptotic approximations” to the
 327 procedure which is actually being employed.

328 There seems to be a prevailing view in statistical ecology that classical
 329 likelihood-based procedures are virtuous because of the availability of simple
 330 formulas and procedures for carrying out inference, such as calculating stan-
 331 dard errors, doing model selection by AIC, and assessing goodness-of-fit. In
 332 large samples, this may be an important practical benefit, but the theoretical
 333 validity of these procedures cannot be asserted in most situations involving small
 334 samples. This is not a minor issue because it is typical in many wildlife sam-
 335 pling problems - especially in surveys of carnivores or rare/endangered species
 336 - to wind up with a small, sometimes extremely small, data set. For example,
 337 a recent paper on the fossa (*Cryptoprocta ferox*), an endangered carnivore in
 338 Madagascar, estimated an adult density of 0.18 adults / km sq based on 20 ani-
 339 mals captured over 3 years (Hawkins and Racey, 2005). A similar paper on the
 340 endangered southern river otter (*Lontra provocax*) estimated a density of 0.25
 341 animals per river km based on 12 individuals captured over 3 years (Sepúlveda

et al., 2007). Gardner et al. (2010) analyzed data from a study of the Pampas cat, a species for which very little is known, wherein only 22 individual cats were captured during the two year period. Trolle and Kéry (2005) reported only 9 individual ocelots captured and Jackson et al. (2006) captured 6 individual snow leopards using camera trapping. Thus, studies of rare and/or secretive carnivores necessarily and flagrantly violate one of Le Cam’s Basic Principles, that of “If you need to use asymptotic arguments, do not forget to let your number of observations tend to infinity.” (Le Cam, 1990).

The biologist thus faces a dilemma with such data. On one hand, these datasets, and the resulting inference, are often criticized as being poor and unreliable. Or, even worse³, “the data set is so small, this is a poor analysis.” On the other hand, such data may be all that is available for species that are extraordinarily important for conservation and management. The Bayesian framework for inference provides a valid, rigorous, and flexible framework that is theoretically justifiable in arbitrary sample sizes. This is not to say that one will obtain precise estimates of density or other parameters, just that your inference is coherent and justifiable from a conceptual and technical statistical point of view. That is, we report the posterior probability $\Pr(D|data)$ which is easily interpretable and just what it is advertised to be and we don’t need to do a simulation study to evaluate how well some approximate $\Pr(D|data)$ deviates from the actual $\Pr(D|data)$ because they are precisely the same quantity.

1.4 Characterizing posterior distributions by MCMC simulation

In practice, it is not really feasible to ever compute the marginal probability distribution $\Pr(y)$, the denominator resulting from application of Bayes’ rule. For decades this impeded the adoption of Bayesian methods by practitioners. Or, the few Bayesian analyses done were based on asymptotic normal approximations to the posterior distribution. While this was useful stuff from a theoretical and technical standpoint and, practically, it allowed people to make the probability statements that they naturally would like to make, it was kind of a bad joke around the Bayesian water-cooler to, on one hand, criticize classical statistics for being, essentially, completely ad hoc in their approach to things but then, on the other hand, have to devise various approximations to what they were trying to characterize. The advent of Markov chain Monte Carlo (MCMC) methods has made it easier to calculate posterior distributions for just about any problem to arbitrary levels of precision.

Broadly speaking, MCMC is a class of methods for drawing random numbers (sampling or simulating) from the target posterior distribution. Thus, even though we might not recognize the posterior as a named distribution or be able to analyze its features analytically, e.g., devise mathematical expressions for the mean and variance, we can use these MCMC methods to obtain a large sample

³Actual quote from a referee

from the posterior and then use that sample to characterize features of the posterior. What we do with the sample depends on our intentions – typically we obtain the mean or median for use as a point estimate, and take a confidence interval based on Monte Carlo estimates of the quantiles. These are estimates, but not like frequentist estimates. Rather, they are Monte Carlo estimates with an associated Monte Carlo error which is largely determined arbitrarily by the analyst. They are not estimates qualified by a sampling distribution as in classical statistics. If we run our MCMC long enough then our reported value of $E[\theta|y]$ or any feature of the posterior distribution is precisely what we say it is. There is no “sampling variation” in the frequentist sense of the word. In summary, the MCMC samples provide a Monte Carlo characterization of *the* posterior distribution.

1.5 What Goes on Under the MCMC Hood

We will develop and apply MCMC methods in some detail for spatial capture-recapture models in chapter ?? . Here we provide a simple illustration of some basic ideas related to the practice of MCMC.

A type of MCMC method relevant to most problems is Gibbs sampling (REF XYZ XYZ), which is based on the idea of iterative simulation from the “full conditional” distributions (also called conditional posterior distributions). The full conditional distribution for an unknown quantity is the conditional distribution of that quantity given every other random variable in the model – the data and all other parameters. For example, for a normal regression model with $y \sim \text{Normal}(\alpha + \beta x, 1)$ then the two full conditionals are, in symbolic terms,

$$[\alpha|y, \beta]$$

and

$$[\beta|y, \alpha].$$

We might use our knowledge of probability to identify these mathematically. In particular, by Bayes’ Rule, $[\alpha|y, \beta] = [y|\alpha, \beta][\alpha|\beta]/[y|\beta]$ and similarly for $[\beta|y, \alpha]$. For example, if we have priors for $[\alpha]$ and $[\beta]$ which are also normal distributions, some algebra reveals that XXXX COPY NOTATION FFROM CH. 6 XXXXX

$$[\alpha|y, \beta] = \text{Normal}(y\text{bar}, \dots \text{weightedvariancehere} \dots).$$

Similarly,

$$[\beta|y, \alpha] \text{isnormal}(\dots \dots \dots)$$

The MCMC algorithm for this model has us simulate in succession, repeatedly, from those two distributions. See ? for more examples of Gibbs sampling for the normal model. A conceptual representation of the MCMC algorithm for this simple model is therefore:

417 Algorithm:

418 0. Initialize α and β

420 Repeat{

421 1. Draw a new value of α from Eq. \ref{xyz}

422 2. Draw a new value of β from Eq. \ref{xyz}

423 }

424 As we just saw for this simple “normal-normal” model it is sometimes possible to specify the full conditional distributions analytically. In general, when certain so-called conjugate prior distributions are chosen, the form of full conditional distributions is similar to that of the observation model. In this normal-normal case, choice of normal priors for the mean parameters is the conjugate prior under the normal model, and thus the full-conditional distributions are also normal. This is convenient because, in such cases, we can simulate directly from them using standard methods (or R functions). But, in practice, we don’t really ever need to know such things because most of the time we can get by using a simple algorithm, called the Metropolis-Hastings (henceforth “MH”) algorithm, to obtain samples from these full conditional distributions without having to recognize them as specific, named, distributions. As we noted above, this gives us enormous freedom in developing models and analyzing them without having to resolve them mathematically because to implement the MH algorithm we need only identify the full conditional distribution up to a constant of proportionality, that being the marginal distribution in the denominator (e.g., $[y|\beta]$ above).

443 1.5.1 Rules for constructing full conditional distributions

444 The basic strategy for constructing full-conditional distributions for devising MCMC algorithms can be reduced conceptually to a couple of basic steps summarized as follows:

445 (step 1) collect all stochastic components of the model;

446 (step 2) Recognize and express the full conditional in question as proportional to the product of all components;

447 (step 3) remove the ones that don’t have the focal parameter in them.

448 (step 4) Do some algebra on the result in order to identify the resulting pdf or pmf.

449 Of the 4 steps, the last of those is the main step that requires quite a bit of statistical experience and intuition because various algebraic tricks can be used to reshape the mess into something noticeable - i.e., a standard, named distribution. But step 4 is not necessary if we decide instead to use the Metropolis-Hastings algorithm as described below.

To illustrate for computing $[\alpha|y, \beta]$ we first apply step 1 and identify the model components as $[y|\alpha, \beta]$, $[\alpha]$ and $[\beta]$. Step 2 has us write $[\alpha|y, \beta] \propto [y|\alpha, \beta][\alpha][\beta]$. We note that $[\beta]$ is not a function of alpha and therefore we delete it to get $[\alpha|y, \beta] \propto [y|\alpha, \beta][\alpha]$. Similarly we get $[\beta|y, \alpha] \propto [y|\alpha, \beta][\beta]$. We can apply step 4 and manipulate these algebraically to arrive at the result or, alternatively, we can sample them indirectly using the Metropolis-Hastings algorithm (see below).

1.5.2 Metropolis-Hastings algorithm

The Metropolis-Hastings algorithm is a completely generic method for sampling from any distribution, say $f(\theta)$. In our applications, $f(\theta)$ will typically be the full conditional distribution for theta. Often, the MH algorithm is used to sample from the full conditional distributions and the resulting synthetic algorithm is called “Metropolis within Gibbs” or similar. Shortly we will actually construct such an algorithm for a simple class of models. The Metropolis-Hastings algorithm generates candidates from some proposal or candidate-generating distribution, that may be conditional on the current value of the parameter, denoted by $h(\theta|\theta^{current})$. Then you accept the proposed value with probability

$$\frac{f(\theta^{cand})h(\theta^{current}|\theta^{cand})}{f(\theta^{current})h(\theta^{cand}|\theta^{current})}$$

this ratio can sometimes be > 1 in which case we set it equal to 1. It is useful to note that $h()$ can be anything at all. Absolutely anything! You can generate candidate values from a *normal*(0,1) distribution, from a uniform(-3455,3455) distribution, or anything of proper support. Note, however, that good choices of $h()$ are those that approximate the posterior distribution. Obviously if $h() = f(\theta|y)$ (i.e., the posterior) then you always accept the draw, and it stands to reason that proposals that are more similar to $f(\theta|y)$ will lead to higher acceptance probabilities. No matter the choice of $h()$, we can evaluate this ratio numerically because the marginal $f(y)$ cancels from both the numerator and denominator. (That is kind of the magic point here that I should emphasize better above.)

A special kind of $h()$ are those that are symmetric, which means that $h(a|b) = h(b|a)$ in which case $h(a|b)$ and $h(b|a)$ just cancel out. A type of symmetric proposal useful in many situations is the so-called “random-walk” proposal distribution where candidate values are drawn from a normal distribution with mean equal to the current value and some standard deviation, say delta which is prescribed by the user. For parameters that have support on the real line, say alpha in our example above, the random walk proposal generator has us generate $\alpha^{*} \sim \text{Normal}(\alpha^{current}, \text{delta})$. If we set delta very small we have a high probability of accepting the proposal and vice versa. In practice, we “tune” delta to achieve a compromise between reasonable mixing of the Markov chains (see below for an example).

Parameters with bounded support: Many models contain parameters that have a bounded support. E.g., variance parameters live on $[0, \infty]$ or similar. In

that case it is sometimes convenient to use a random walk proposal distribution, but just reject parameters that are outside of the parameter space (REF FOR THIS?).

1.6 Practical Bayesian Analysis and MCMC

There are a number of really important practical issues to be considered in any Bayesian analysis and we cover some of these briefly here.

Prior distributions: Bayesian analysis requires that we choose prior distributions for all of the structural parameters of the model (we use the term structural parameter to mean all parameters that aren't customary thought of as latent variables). We will strive to use priors that are meant to express little or no prior information - default or customary "non-informative" or diffuse priors. This will be uniform(a,b) priors for parameters that have a natural bounded support and, for parameters that live on the real line we use either (1) diffuse normal priors; (2) "improper" uniform priors or (3) sometimes even a bounded uniform(a,b) prior if that greatly improves the performance of WinBUGS or other software doing the MCMC for us. In WinBUGS a prior with low "precision" (precision = 1/sigma2) such as normal(0,.01) will typically be used. Of course tau = 0.01 (sigma2 = 100) might be very informative for a regression parameter that has a high variance. Therefore, we recommend that predictor variables *always* be standardized. Clearly there are a lot of choices for ostensibly non-informative priors, and the degree of non-informativeness depends on the parameterization. For example, a natural non-informative prior for the intercept of a logistic regression

$$\text{logit}(p[i]) = a + b * x[i]$$

Would be $[a] = \text{const}$ which is the same as saying $a \sim \text{Unif}(\infty, \text{infity})$ or the standard improper "locally uniform" prior distribution. However, we might also use a prior on the parameter $p0 = \text{expit}(a)$, which is $Pr(y = 1)$ for the value $x = 0$. Since $p0$ is a probability we might use $p0 \sim \text{Unif}(0, 1)$. These two priors can affect results (see Chapter 3.XYZ), yet they are both sensible "non-informative" priors. Choice of priors and parameterization is very much problem-specific and often largely subjective. Moreover, it also affects the behavior of MCMC algorithms and therefore the analyst needs to pay some attention to these issues and possibly try different things out. [we should point to some standard refs on this stuff].

Once we have carried-out an analysis by MCMC, there are many other practical issues that we have to confront. One of the most important is "Have the chains converged?" Most MCMC algorithms only guarantee that, eventually, the samples being generated will be from the target posterior distribution. So-called "convergence" of the Markov chain is achieved when that happens. Typically a period of transience is observed in the early part of the MCMC algorithm, and this is usually discarded as the "burn-in" period.

The quick diagnostic to whether convergence has been achieved is that your Markov chains look "grassy" - see Figure XXX below - then you're probably all

done. Another way to check convergence is to update the parameters some more and see if the posterior changes. It is good to confirm convergence using the Rhat statistic (Brooks Gelman Rubin statistic (Gelman et al., 1996)) which should be close to 1. In practice, 1.2 is probably good enough. For some really complex models 1.3 or 1.4 might be good enough. For some models you can't actually realize a low R-hat. E.g., if the posterior is a discrete mixture of distributions then I think you will always be misled into thinking that your Markov chains have not converged when in fact the chains are just jumping back and forth in the posterior state-space. Another situation is when one of the parameters is on the boundary of the parameter space which might appear to be very poor mixing. This kind of stuff is normally ok and you need to think really hard about the context of the model and the problem before you conclude that your MCMC algorithm is ill-behaved or not.

Some models exhibit "poor mixing" of the Markov chains or what people might also call "slow convergence" which is a term we would disagree with because the samples might well be from the posterior (i.e., the Markov chains have converged to the proper stationary distribution) but simply mix around the posterior rather slowly. Anyway, poor mixing can happen for a huge number of reasons - when parameters are highly correlated (even confounded), or barely identified from the data, or the algorithms are very terrible and probably many other reasons. Slow mixing equates to high autocorrelation in the Markov chain - the successive draws are highly correlated, and thus we need to run the MCMC algorithm much longer to get an effective sample size that is sufficient for estimation - or to reduce the MC error to a tolerable level. A strategy often used to reduce autocorrelation is "thinning" - i.e., keep every m^{th} value of the Markov chain output. However, thinning is necessarily inefficient from the stand point of inference - you can always get more precise posterior estimates by using all of the MCMC output regardless of the level of autocorrelation (MacEachern and Berliner, 1994). Practical considerations might necessitate thinning, even though it is statistically inefficient. For example, in models with many parameters or other unknowns being tabulated, the output files might be enormous and unwieldy to work with. In such cases, thinning is perfectly reasonable. In many cases, how well the Markov chains mix is strongly influenced by parameterization, standardization of covariates, and the prior distributions being used. Some things work better than others, and the investigator should experiment with different settings and try not to become bewildered when things don't work out perfectly. MCMC is an art, and a science.

The next question: Is the posterior sample large enough? Never report MCMC results to more than 2 decimal places - because they will always be different! Look at the MC error which is printed by default in *BUGS summaries. You want that to be smallish relative to the magnitude of the parameter. I'm usually content with 1% but if you're uncomfortable with monte carlo error, you should run your MCMC algorithm as long as it takes. Note that MC error in summaries of the posterior is not the same as having an "approximate" solution in a standard likelihood analysis or similar. The approximate SE in likelihood inference is actually wrong in its actual value.... XYZ.

587 1.6.1 Bayesian confidence intervals

588 The 95% Bayesian interval based on percentiles of the posterior is not a unique
 589 interval - there are many of them - and the so-called “highest posterior density”
 590 (HPD) interval is the narrowest interval. We might compute that frequently
 591 because it is easy to do with an integer parameter which N is (See the next
 592 chapter). The 95p% HPD is not often exactly 95% but usually slightly more
 593 conservative than nominal because it is the narrowest interval that contains at
 594 least 95% of the posterior mass.

595 1.6.2 Estimating functions of parameters

596 A benefit of analysis by MCMC is that we can seamlessly estimate functions
 597 of parameters by simply tabulating the desired function of the simulated pos-
 598 terior draws. For example, if θ is the parameter of interest and let $\theta^{(i)}$ for
 599 $i = 1, 2, \dots, M$ be the posterior samples of θ . Let $\eta = \exp(\theta)$, then a posterior
 600 sample of η can be obtained simply by computing $\exp(\theta^{(i)})$ for $i = 1, 2, \dots, M$.
 601 We give an example in Section XXXX below.

602 1.7 Bayesian Analysis using WinBUGS

603 We won’t be too concerned with devising our own MCMC algorithms although
 604 we will do that one or two times for fun. More often, we will rely on the
 605 freely available software package WinBUGS or other BUGS engines for doing
 606 this. Further, we will execute WinBUGS from within R using the R2WinBUGS
 607 package. WinBUGS is an MCMC black box that takes a pseudo-code descrip-
 608 tion of all of the relevant stochastic and deterministic elements of a model and
 609 generates an MCMC algorithm for that model. But you never get to see the al-
 610 gorithm. Instead, WinBUGS will run the algorithm and just return the Markov
 611 chain output - the posterior samples of model parameters.

612 The great thing about WinBUGS is that it forces you to become intimate
 613 with your statistical model - you have to write each element of the model down,
 614 admit (explicitly) all of the various assumptions, understand what the actual
 615 probability assumptions are and how data relate to latent variables and data and
 616 latent variables relate to parameters, and how parameters relate to one another.
 617 While we will use WinBUGS almost exclusively here, there are many BUGS like
 618 packages now, including JAGS, OpenBUGS, PyMC and others. Later (chap-
 619 ter MCMC XYZ) we will demonstrate a model or two in JAGS. OpenBUGS
 620 is the current active development tree of the “BUGS” language. See (Kéry
 621 (2010); chapters XXXX) and (Kery and Schaub (2011), Appendix XYZ) for the
 622 lowdown on problems/issues with using WinBUGS. That book should also be
 623 consulted for a more comprehensive introduction to using WinBUGS. In this
 624 example, we’re going to accelerate pretty fast.

625 We provide a brief introductory example of a normal regression model using
 626 a small simulated data set. The following commands are executed from within

your R workspace, the command line being indicated by “`>`”. First, simulate a covariate x and observations y having prescribed intercept, slope and variance:

```

629 > x<-rnorm(10)
630 > mu<- -3.2+ 1.5*x
631 > y<-rnorm(10,mu,sd=4)

```

The WinBUGS model specification for a normal regression model is written within R as a character string input to the command `cat()` and then dumped to a text file named “normal.txt” (alternatively, you can write the model specifications directly within a text file and save it in your current working directory):

```

636 > cat("
637 model {
638   for (i in 1:10){
639     y[i]~dnorm(mu[i],tau)           # the "likelihood"
640     mu[i]<- beta0 + beta1*x[i]      # the linear predictor
641   }
642   beta0~dnorm(0,.01)               # prior distribution
643   beta1~dnorm(0,.01)
644   sigma~dunif(0,100)
645   tau<-1/(sigma*sigma)             # tau is a derived parameter
646 }
647 ",file="normal.txt")

```

Remarks:

1. WinBUGS parameterizes the normal in terms of the mean and inverse-variance, called the precision. Thus, `dnorm(0,.01)` implies a variance of 100.
2. We typically use diffuse normal priors for mean parameters, `beta0` and `beta1` in this case, but sometimes we might use uniform priors with suitable bounds `-B` and `+B`.
3. We typically use a uniform `[0,B]` prior on standard deviation parameters (Gelman XXX 2006). But sometimes we might use a gamma prior on the precision parameter `tau`.
4. In a WinBUGS model file, every single element has to be either data which will be input (see below), a random variable which must have a probability distribution associated with it, using the “`~`”, or it has to be a derived parameter connected to variables and data using “`-`”.

To fit the model, we execute these commands:

```

663 > library("R2WinBUGS")      # "attach" the R2WinBUGS library
664 > data <- list ( "y","x")
665 > inits <- function()

```

```

666 list ( beta1=rnorm(1),beta0=rnorm(1),sigma=runif(1,0,2) )
667 > parameters <- c("beta0","beta1","sigma","tau")
668 > out<-bugs (data, inits, parameters, "normal.txt", n.thin=2, n.chains=2, n.burnin=2000)

```

669 To fit the model, we execute these commands:

```

670 > library("R2WinBUGS")      # "attach" the R2WinBUGS library
671 > data <- list ( "y","x")
672 > inits <- function()
673   list ( beta1=rnorm(1),beta0=rnorm(1),sigma=runif(1,0,2) )
674 > parameters <- c("beta0","beta1","sigma","tau")
675 > out<-bugs (data, inits, parameters, "normal.txt", n.thin=2, n.chains=2, n.burnin=2000)

```

676 **Explanation:** We created an R list object called “data” which are the
 677 things we have to send to WinBUGS. In the example above, the data consist
 678 of two objects which exist as “y” and “x” in the R workspace and also in the
 679 WinBUGS model definition. People tend to ask “how should my data be format-
 680 ted?” That depends on how you describe the WinBUGS model and you should
 681 read your data in as a .csv file or some other format and manipulated it within
 682 R to get into the desired format. There is a non-unique way to describe any
 683 particular model and so you have some flexibility. We talk about data format
 684 further in the context of capture-recapture models and SCR models in chapters
 685 3 and 4, and later. We also have to create an R function that produces a list
 686 of starting values “inits” that get sent to WinBUGS. In general, starting values
 687 are optional but we recommend to always provide reasonable starting values of
 688 structural parameters, but not necessarily random effects(although the latter
 689 will sometimes need to be given to keep WinBUGS from crashing). Finally, we
 690 identify the names of the parameters (labeled correspondingly in the WinBUGS
 691 model specification) that we want WinBUGS to save the MCMC output for. In
 692 the above example, we are telling WinBUGS to “monitor” beta0, beta1, sigma
 693 and tau. WinBUGS is executed using the R command “bugs”. Note that the
 694 previously created objects defining data, initial values and parameters to mon-
 695 itor are passed to this function. In addition, various other things are declared:
 696 The number of chains, the thinning rate, the number of burnin iterations and
 697 the total number of iterations. We set “debug=TRUE” if we want the Win-
 698 BUGS GUI to stay open (useful for analyzing MCMC output and looking at
 699 the WinBUGS error log). Also, we set working.dir=getwd() so that WinBUGS
 700 output files and the log file are saved in the current R working directory.

701 You should execute all of the commands given above and then look at the
 702 resulting output. Kill the WinBUGS GUI and the data will be read back into
 703 R. We don’t want to give instructions on how to navigate and use the GUI - see
 704 REF (XYZ) for that. The object “out” prints important summaries by default
 705 (this is slightly edited):

```

706 > print(out,digits=2)
707 Inference for Bugs model at "normal.txt", fit using WinBUGS,
708 2 chains, each with 6000 iterations (first 2000 discarded), n.thin = 2

```

```

709  n.sims = 4000 iterations saved
710      mean   sd 2.5%   25%   50%   75% 97.5% Rhat n.eff
711  beta0    -2.43 1.84 -6.21 -3.50 -2.42 -1.34 1.27   1 4000
712  beta1     2.62 1.54 -0.42  1.68  2.62  3.57  5.67   1 4000
713  sigma     5.29 1.66  3.11  4.14  4.95  6.05  9.39   1 4000
714  tau       0.05 0.02  0.01  0.03  0.04  0.06  0.10   1 4000
715  deviance 59.85 3.24 56.18 57.47 59.00 61.37 68.32   1 840
716
717  For each parameter, n.eff is a crude measure of effective sample size,
718  and Rhat is the potential scale reduction factor (at convergence, Rhat=1).
719
720  DIC info (using the rule, pD = Dbar-Dhat)
721  pD = 2.6 and DIC = 62.4

```

722 **Remarks:** (1) convergence is assessed using the \hat{R} statistic - which we
723 will write “Rhat”. A value of Rhat near 1 indicates convergence. Posterior
724 summaries are given. (2) DIC is the “deviance information criterion” (REF
725 XXXX; see below XYZ) which some people use in a manner similar to AIC
726 although it is recognized to have some problems in hierarchical models (XYZ
727 Biometrics ref XYZ).

728 **Inference about functions of model parameters:** Using the MCMC
729 draws for a given model we can easily obtain the posterior distribution of any
730 function of model parameters. We showed this by providing the posterior of
731 “tau” when we used “sigma” to parameterize the model above. As another
732 example, suppose that the normal regression model above had a quadratic re-
733 sponse function of the form

$$E[y[i]] = \beta_0 + \beta_1 * x[i] + \beta_2 * x[i] * x[i]$$

734 Then the optimum response can be found by setting the derivative of this func-
735 tion to 0 and solving for x . We find that $df/dx = \beta_1 + 2 * \beta_2 * x = 0$
736 yields that $x_{opt} = -\beta_1/(2 * \beta_2)$. We can just take our posterior draws for β_1
737 and β_2 and obtain a posterior sample of x_{opt} using those values. As an exer-
738 cise, take the normal model above and simulate a quadratic response and then
739 describe the posterior distribution of x_{opt} .

740 1.8 Model Checking and Selection

741 In general terms model checking - or assessing the adequacy of the model - and
742 model selection are quite thorny issues and, despite contrary and commonly held
743 belief among practitioners, there are not really definitive, general solutions to
744 either problem. We’re against dogma on these issues and think people need to be
745 open-minded about such things and recognize that models can be useful whether
746 or not they pass certain statistical tests. Some models are intrinsically better
747 than others because they make more biological sense or foster understanding
748 or achieve some objective that a bootstrap goodness-of-fit test can’t decide for

you. In the context of Bayesian model checking and selection see Kéry (2010);
chapter XYZ, and Link and Barker (2009); chapter XYZ.

1.8.1 Goodness-of-fit

Goodness-of-fit testing is an important element of any analysis because in a sense our model represents a general set of hypotheses about the ecological and observation processes that generated our data. Thus, if our model “fits” in some statistical or scientific sense, then we believe it to be consistent with the hypotheses that went into the model. More formally, we would conclude that the data are *not inconsistent* with the hypotheses. If we have enough data, then of course we will reject any set of statistical hypotheses. Unfortunately, conducting goodness-of-fit tests is not always so easy to do. Moreover, it is never really easy (or especially convenient) to decide if your goodness-of-fit test is worth anything. It might have 0 power! Despite these difficulties, we will often try to conjure something up that gets the job done.

Even though we think evaluation of fit is important, we also believe that models can be useful irrespective of whether they fit (as we noted above, with enough data, no model will fit, and some contributing factors to lack-of-fit can be minor or irrelevant to the intended use of the model). As a final point, we can always make a model fit by making the model extremely complex. It seems to us that simple models that you can understand should usually be preferred even if they don’t fit. Yet the tension is there to get fitting models which comes naturally at the expense of models that can be interpreted and studied and used.

To evaluate goodness-of-fit in Bayesian analyses, we will most often use the Bayesian p-value (Gelman XXYZZ). The basic idea is to define a fit statistic and compare the posterior distribution of that statistic to the posterior predictive distribution of that statistic for hypothetical perfect data sets for which the model is correct. For example, with count frequency data, a standard measure of fit is the sum of squares of the “Pearson residuals”,

$$D[i] = (y[i] - E[y[i]])^2 / \text{Var}[y[i]]$$

The fit statistic based on the squared residuals is

$$FIT = \sum_i D[i]^2$$

which can be computed at each iteration of a MCMC algorithm given the current values of parameters that determine the mean and variance of the response distribution. The equivalent statistic is computed for a “new” data set, simulated using the current parameter values. The Bayesian p-value is simply the posterior probability $Pr(FIT > FIT_{new})$ which should be close to 0.50 for a good model. In practice we judge “close to 0.50” as being “not too close to 0 or 1” and, as always, closeness is somewhat subjective. We’re happy with anything $> .1$ and $< .9$ but might settle for $> .05$ and < 0.95 . In summary, the Bayesian

p-value seems like a bootstrap idea, is easy to compute, and widely used as a result.

Sometimes a more useful fit statistic is the Freeman-Tukey statistic, in which

$$D(x, \theta) = \sum_j (\sqrt{x_j} - \text{sqr}t(e_j))^2$$

(Brooks et al., 2000), where x_j is the observed value of observation j and e_j its expected value. In contrast to a chi-square discrepancy, the Freeman-Tukey statistic removes the need to pool cells with small expected values.

1.8.2 Model Selection

For model selection we typically use three different methods: First is, let's say, common sense. If a parameter has posterior mass concentrated away from 0 then it seems like it should be regarded as important - that is, it is "significant." This approach seems to have fallen out of favor with all of the interest over the last 10 or 15 years on model selection in ecology. It seems reasonable to us.

For regression problems we use the factor weighting idea which is to introduce a set of binary variables $w(k)$ for variable k , and express the model as, e.g., for a single covariate model:

$$E[y[i]] = a + w * b * x[i]$$

where w is given a Bernoulli prior distribution with some prescribed probability. E.g., $w \sim \text{Bern}(0.50)$ to provide a prior probability of 0.50 that variable "x" should be an element of the linear predictor. The posterior probability of the event $w = 1$ is a gauge of the importance of the variable $x[i]$. i.e., high values of $\text{Pr}(w = 1)$ indicate stronger evidence....close to 0 means not so important, etc... This idea seems to be due to Kuo and Mallick (XXX)⁴ and see Royle and Dorazio (2008); ch XX for an example in the context of logistic regression. It seems to even work sometimes with fairly complex hierarchical models of a certain form. E.g., Royle (2008) applied it to a random effects model where w multiplied the random effect. WinBUGS can be very sensitive and temperamental to things but sometimes it does things that appear to be quite remarkable. The problem with this approach is that its effectiveness and results will typically be highly sensitive to the prior distribution on the structural parameters (e.g., see Royle and Dorazio (2008) table XYZ). The reason for this is obvious: If $w = 0$ for the current iteration of the MCMC algorithm, so that "b" is sampled from the prior distribution, and the prior distribution is very diffuse, then extreme values of "b" are likely. When the current value of "b" is far away from the mass of the posterior when $w = 1$, then the Markov chain may only jump from $w = 0$ to $w = 1$ infrequently. One seemingly reasonable solution to this problem (Aitken XYZ) is to fit the full model to obtain posterior distributions for all parameters, and then use those as prior distributions in a "model selection" run

⁴Is this also what people call Zellner's G-priors?

of the MCMC algorithm. This seems preferable to an arbitrary restriction of the prior support to improve the performance of the MCMC algorithm.

A third method that we like to fall-back on is subject-matter context. It seems that there are some situations where one should not have to do model selection because it is necessitated by the specific situation at hand. SCR models are such an example. We will see that “spatial location” of individuals is an element of the model. The simpler, reduced, model is an ordinary capture-recapture model (i.e., next chapter), but it seems silly to think about actually using the reduced model even if we could concoct some statistical test to refute the more complex model. Other examples are when effort, area or sample rate is a covariate. One might prefer to have such things in models regardless of whether or not they pass some statistical litmus test (yet you can always find referees to argue for pedantic procedure over thinking).

Many problems can be approached using one of these methods but there are also broad classes of problems that can’t and, for those, you’re out of luck. In later chapters we will address model selection in specific contexts and we hope those will prove useful.

1.9 Poisson GLMs

The Poisson GLM (also known as “Poisson regression”) is probably the most relevant and important class of models in all of ecology. The basic model assumes observations $y(i); i = 1, 2, \dots, n$ follow a Poisson distribution with mean λ which we write

$$y(i) \sim \text{Poisson}(\lambda)$$

Commonly $y(i)$ is a count of animals or plants at some point in space and λ might depend on i . For example, i might index point count locations in a forest, BBS route centers, or sample quadrats, or similar. If covariates are available it is typical to model them as linear effects on the log mean. If $x(i)$ is some measured covariate associated with observation i . Then,

$$\log(x(i)) = \alpha + \beta * x(i)$$

While we only specify the mean of the Poisson model directly, the Poisson model (and all GLMs) has a “built-in” variance which is directly related to the mean. In this case, $\text{Var}(y) = E(y) = \lambda$. Thus the model accommodates a linear increase in variance with the mean. Another extremely useful feature of the Poisson model is the property of “compound additivity”. If $y(1)$ and $y(2)$ are Poisson random variables with means $\lambda[1]$ and $\lambda[2]$, then $y(1) + y(2)$ is Poisson with mean $(\lambda[1] + \lambda[2])$. Thus, if the observations can be viewed as an aggregate of counts over some finer scale, then the mean aggregates in a corresponding manner. Multinomial random variables have a direct relationship to Poisson random variables. If $y(1)$ and $y(2)$ are *iid* Poisson then, conditional on their total $T = y(1) + y(2)$, they have a multinomial distribution with sample size T and cell probabilities $\lambda[1]/(\lambda[1] + \lambda[2])$ and $\lambda[2]/(\lambda[1] + \lambda[2])$. These are some of the reasons the Poisson distribution is extremely useful in ecology.

1.9.1 Example: Breeding Bird Survey Data

As an example we consider a classical situation in ecology where counts of an organism are made at a collection of spatial locations. In this particular example, we have mourning dove counts made along North American Breeding Bird Survey (BBS) routes in Pennsylvania, USA. A route consists of 50 stops separated by 0.5 mile. For the purposes here we are defining $y[i]$ = route total count and the sample location will be marked by the center point of the BBS route. The survey is run annually and the data set we have is 1966-1998. BBS data can be obtained online at <http://...xyz.xyz.xyz>. We will make use of the whole data set shortly but for now we're going to focus on a specific year of counts - 1990 - for no particular reason. For 1990 there were 77 active routes. We have the data stored in a .csv file where rows index the unique route, column 1 is the route ID, columns 2-3 are the route coordinates (longitude/latitude), column 4 is a habitat covariate "forest cover" (standardized, see below) and the remaining columns are the yearly counts. Years for which a route was not run are coded as "NA" in the data matrix. We imagine that this will be a typical format for many ecological studies, perhaps with more columns representing covariates. To read in the data and display the first few elements of this matrix, do this:

```
> a<-read.csv("pa-bbsdovedata-all.csv")
> data[1:2,1:6]
      X      lon      lat      habitat X66 X67
1 72002 -80.445 41.501 -0.3871372   NA   24
2 72003 -80.347 41.214 -1.0171629   NA   NA
```

It is useful to display the pattern in counts. For that we use a spatial dot plot - where we plot the coordinates of the observations and mark the color of the plotting symbol based on the magnitude of the count. We have a special plotting function for that which is called `spatial.plot()` and it is available with the supplemental materials. Actually, what we want to do here is plot the log-count (+1 of course!) which displays a notable pattern that could be related to something. We can ponder the potential effects that might lead to dove counts being high....Corn fields, telephone wires, barn roofs along with misidentification of pigeons, these could all correlated reasonably well with these counts for all we know. Unfortunately we don't have any of that information.

We do have a measure of forest cover in the vicinity of each point which is contained in the data set ("habitat"). This was derived from a larger GIS coverage of the state (provided in the data file "pahabdata") which can be plotted using the `spatial.plot` function using the following commands

```
> map('state',regions="penn",lwd=2)
> spatial.plot(pahabdata[,2:3],pahabdata[, "dfor"],cx=2)
> map('state',regions="penn",lwd=2,add=TRUE)
```

We see a prominent pattern that indicates high forest coverage in the central part of the state and low forest cover in the SE. Inspecting the previous figure

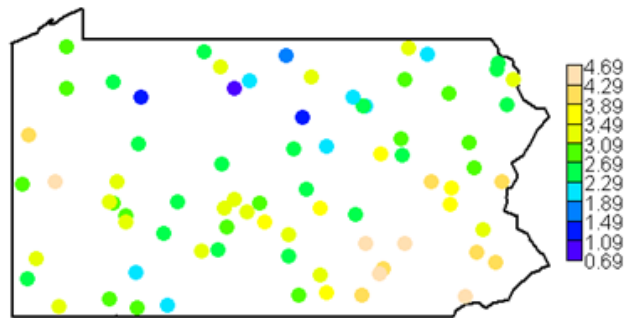


Figure 1.1: Needs a caption

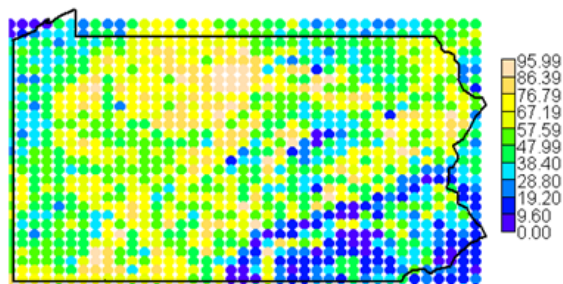


Figure 1.2: Needs a caption

906 of log-counts suggests a relationship between counts and forest cover which is
 907 not surprising.

908 1.9.2 Doing it in WinBUGS

909 Here we demonstrate how to fit a Poisson GLM in WinBUGS using the covariate
 910 $x(i)$ = forest cover. It is advisable that $x(i)$ be standardized in most cases as
 911 this will improve mixing of the Markov chains. Recall that the data we have
 912 stored include a standardized covariate (forest cover) and so we don't have to
 913 worry about that here. To read the BBS data into R and get things set up for
 914 WinBUGS we issue the following commands:

```
915 data<-read.csv("pa-bbsdovedata-all.csv")
916 y<-data[,29] # pick out 1990
917 notna<-!is.na(y)
918 y<-y[notna]
919 habitat<-data[notna,4]
920 library("R2WinBUGS")
921 data <- list ( "y","M","habitat")
```

922 Now we write out the Poisson model specification in WinBUGS pseudo-code,
 923 provide initial values, identify parameters to be monitored and then execute
 924 WinBUGS:

```
925 cat("
926 model {
927     for (i in 1:M){
928         y[i]~dpois(lam[i])
929         log(lam[i])<- beta0+beta1*habitat[i]
930     }
931     beta0~dunif(-5,5)
932     beta1~dunif(-5,5)
933 }
934 ",file="PoissonGLM.txt")
935
936 inits <- function() list ( beta0=rnorm(1),beta1=rnorm(1))
937 parameters <- c("beta0","beta1")
938 out<-bugs (data, inits, parameters, "PoissonGLM.txt", n.thin=2, n.chains=2, n.burnin=2000,n.iter=
```

939 **Remarks:** (1) Note the close correspondence in how the model is specified
 940 here compared with the normal regression model previously. As an exercise you
 941 should discuss the specific differences between the BUGS model specifications
 942 for the normal and Poisson models.

```
943 > print(out,digits=3)
944 Inference for Bugs model at
945 ‘‘PoissonGLM.txt’’, fit using WinBUGS,
946 2 chains, each with 4000 iterations (first 1000 discarded), n.thin = 2
```

```

947 n.sims = 3000 iterations saved
948      mean      sd      2.5%      25%      50%      75%      97.5%  Rhat  n.eff
949 beta0      3.151  0.025      3.102      3.135      3.151      3.168      3.199 1.001  2300
950 beta1     -0.498  0.021     -0.539     -0.512     -0.498     -0.484     -0.457 1.001  3000
951 fit      869.930 19.856    835.500    855.700    868.600    881.900    913.602 1.002  1600
952 fitnew     76.709 12.519     54.098     68.107     76.215     84.510    102.602 1.001  3000
953 deviance 1116.605  2.014   1115.000   1115.000   1116.000   1117.000   1122.000
954 1.001  3000

```

955 We might wonder whether this model provides an adequate fit to our data.
956 To evaluate that, we used a Bayesian p-value analysis with fit statistic based
957 on the Freeman-Tukey residual by replacing the model specification above with
958 this:

```

959 cat("
960 model {
961   for (i in 1:M){
962     y[i]~dpois(lam[i])
963     log(lam[i])<- beta0+beta1*habitat[i]
964     d[i]<- pow(pow(y[i],0.5)-pow(lam[i],0.5),2) #
965
966     ynew[i]~dpois(lam[i])
967     dnew[i]<-pow( pow(ynew[i],0.5)-pow(lam[i],0.5),2)
968
969   }
970   fit<-sum(d[])
971   fitnew<-sum(dnew[])
972   beta0~dunif(-5,5)
973   beta1~dunif(-5,5)
974 }
975
976
977 ",file="PoissonGLM.txt")

```

978 The Bayesian p-value is the proportion of times $fitnew > fit$ which, for this
979 data set, is 0, which was 1.0 in this case (calculation omitted). This suggests
980 that the basic Poisson model does not fit well.

981 1.9.3 Constructing your own MCMC algorithm

982 It will be helpful for people to suffer through a couple examples building a
983 custom MCMC algorithm. So, here, we build a basic one for the Poisson regres-
984 sion model using a Metropolis-within-Gibbs approach. First, we will assume
985 that the two parameters have diffuse normal priors, say $[\alpha] = norm(0, 100)$ and
986 $[\beta] = norm(0, 100)$. We need to collect the relevant elements of the model which
987 are the likelihood $[y|\alpha, \beta] = prod_i [y[i]|\alpha\beta]$ which is, mathematically, the prod-
988 uct of the Poisson pmf evaluated at $y[i]$, given particular values of β_0 and β_1 .

989 The priors are $[\alpha]$ and $[\beta]$. We identify the full conditionals which are $[\alpha|\beta, y]$
 990 and $[\beta|\alpha, y]$. We use the all-purpose rule for constructing full conditionals to
 991 discover that:

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

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

993 Remember we could replace the “propto” with “equals” if we simply put $[y|\beta]$
 994 or $[y|\alpha]$ in the denominator. But, in general, $[y|\alpha]$ or $[y|\beta]$ will be quite a
 995 pain to compute and, more importantly, it is a constant as far as the operative
 996 parameter (beta or alpha, respectively) goes so we can just as well ignore it
 997 because, recall, the MH acceptance probability will be the ratio of the full-
 998 conditional evaluated at a candidate draw to that evaluated at the current
 999 draw. So, the denominator required to change α to $=$ winds up canceling from
 1000 the MH acceptance probability. Here we will use the random walk candidate
 1001 generator. The “Metropolis within Gibbs” algorithm for a Poisson regression is
 1002 remarkably simple:

1003 I would break this code up into more lines and have objects called “prior” and “prior.candidate”

1004

1005 You could also mention that this is a random walk M-H. It would help lots of people out to see a

1006

```

1007 # put random number seed here
1008 out<-matrix(NA,nrow=1000,ncol=2)    # matrix to store the output
1009 beta0<- -1                          # starting values
1010 beta1<- -.8
1011
1012 # begin the MCMC loop ; do 1000 iterations
1013 for(i in 1:1000){
1014
1015   # update the beta0 parameter
1016   lik.curr<- sum(log(dpois(y,exp(beta0+beta1*habitat))))
1017   prior.curr<- log(dnorm(beta0,0,100))
1018   beta0c<-rnorm(1,beta0,.25)        # generate candidate
1019   lik.cand<- sum(log(dpois(y,exp(beta0c+beta1*habitat))))
1020   prior.cand<- log(dnorm(beta0c,0,100))
1021   if(runif(1)< exp(lik.cand+prior.cand-lik.curr-prior.curr)) beta0<-beta0c
1022
1023   # update the beta1 parameter
1024   lik.curr<- sum(log(dpois(y,exp(beta0+beta1*habitat))))
1025   prior.curr<- log(dnorm(beta1,0,100))
1026   beta1c<-rnorm(1,beta1,.25)
1027   lik.cand<- sum(log(dpois(y,exp(beta0+beta1c*habitat))))
1028   prior.cand<- log(dnorm(beta1c,0,100))
1029   if(runif(1)< exp(lik.cand+prior.cand-lik.curr-prior.curr)) beta1<-beta1c
1030   out[i,<-c(beta0,beta1)             # save the current values
1031 }
```

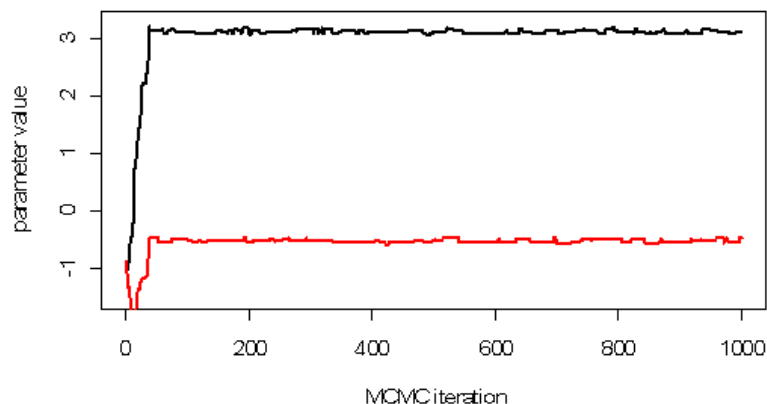


Figure 1.3: Needs a caption

Look at the output (beta0 in red, beta1 in black). You might not like the appearance of this output too much but a couple of things are evident: The Markov chains clearly stabilize - “converge” – after about 100 iterations. They also appear to mix very slowly, although this is not so clear given the scale of the y-axis.

We decreased the variance for candidate generating distribution and re-ran the MCMC algorithm producing the history plots below. We see that the burn-in takes longer but it seems to mix better.

Fig. XYZ shows a longer MCMC run (10,000 total iterations) for beta1 based on discarding the first 400 samples as burn-in. The “grassy” look of the MCMC history is diagnostic of Markov chains that are well-mixing.

Remarks: We used a specific set of starting values for these simulations. It should be clear that starting values closer to the mass of the posterior distribution might cause burn-in to occur faster. As an exercise, evaluate that. (2) Clearly the influence of the proposal variance term is important. Small values lead to much better mixing but it should be noted that values that are too small will lead to very slow mixing. We saw that values that were too large tended to get the parameters stuck in one spot. This suggests there is an optimal value of the Metropolis-Hastings tuning parameter⁵. As an exercise you should find that optimal value. (3) For the flat normal prior distributions here we could leave the prior contribution out of the full conditional evaluation since it is “locally constant”. Note also that we have used a different prior than in our WinBUGS model specification. As an exercise, evaluate whether this seems to affect the

⁵Defined previously?

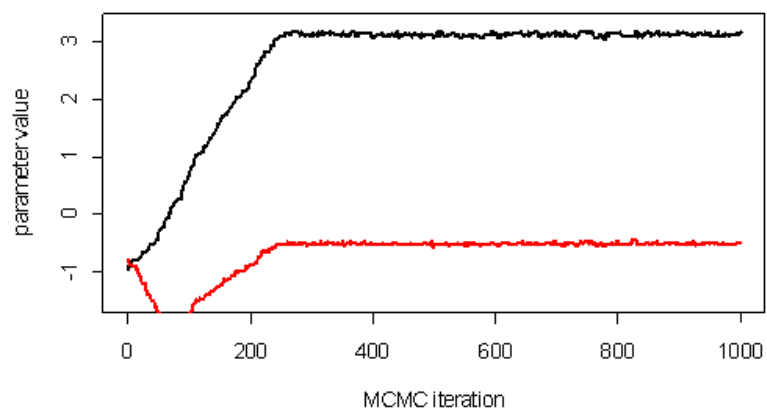


Figure 1.4: Needs a caption

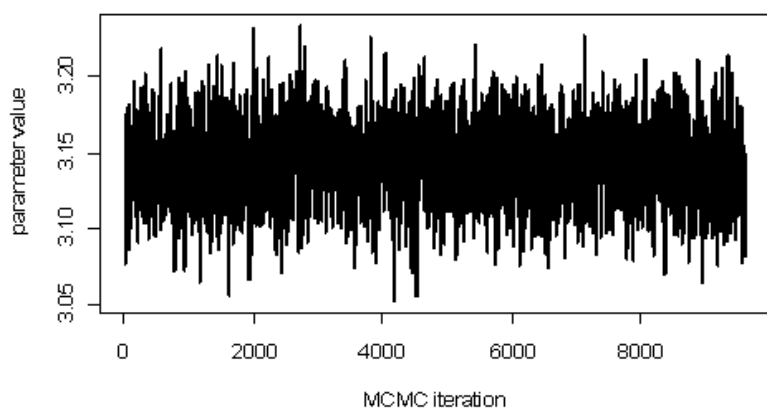


Figure 1.5: Needs a caption

1055 result.

1056 1.10 Poisson GLM with Random Effects

1057 What we will be doing in most of this book is dealing with random effects
1058 in GLM-like models - models that are usually referred to as generalized linear
1059 mixed models (GLMMs).

1060 **The Log-Normal mixture:** The classical situation involves a GLM with
1061 a normally distributed random effect. The linear predictor of the Poisson model
1062 is extended simply by adding a noise term, say:

$$\log(\lambda(i)) = \alpha + \beta * x(i) + \eta[i]$$

1063 where $\eta[i]$ *normal*(0, σ^2). A natural alternative is to have $\exp(\eta[i]) / \sim \gamma(a, b)$
1064 which would correspond to a negative binomial kind of over-dispersion whereas
1065 the normal noise has a different mean/variance relationship (the interested
1066 reader should work that out). Choosing between such possibilities is not a
1067 topic we will get into here because it doesn't seem possible to provide general
1068 guidance on it. Anyhow, it is really amazingly simple to express this model in
1069 WinBUGS and have WinBUGS draw samples from the posterior distribution
1070 using the following code for the BBS dove counts:

```
1071 data<-read.csv("pa-bbsdovedata-all.csv")
1072 locs<-data[,2:3]
1073 habitat<-data[,4]
1074 y<-data[,29]
1075 notna<-!is.na(y) # to remove missing values
1076 y<-y[notna]
1077 locs<-locs[notna,]
1078 habitat<-habitat[notna]
1079 M<-length(y)
1080
1081 cat("
1082 model {
1083     for (i in 1:M){
1084         y[i]~dpois(lam[i])
1085         log(lam[i])<- beta0+beta1*habitat[i] + eta[i]
1086         eta[i] ~ dnorm(0,tau)
1087     }
1088     beta0~dunif(-5,5)
1089     beta1~dunif(-5,5)
1090     sigma~dunif(0,10)
1091     tau<-1/(sigma*sigma)
1092 }
```

1093 I have removed the final several R commands which package up the data and
1094 execute WinBUGS as those commands are largely redundant with the previous

demo. The summary results are:

```

> print(out,digits=3)
Inference for Bugs model at "model.txt", fit using WinBUGS,
  2 chains, each with 5000 iterations (first 1000 discarded), n.thin = 2
 n.sims = 4000 iterations saved
      mean      sd    2.5%    25%    50%    75%   97.5%  Rhat  n.eff
beta0    2.967  0.076   2.817   2.915   2.969   3.020   3.111 1.006   430
beta1   -0.518  0.073  -0.657  -0.566  -0.517  -0.470  -0.374 1.008  4000
sigma    0.598  0.059   0.491   0.556   0.594   0.634   0.725 1.004   640
tau      2.883  0.569   1.904   2.489   2.836   3.233   4.149 1.004   640
fit      19.885  3.190  14.119  17.670  19.705  21.902  26.610 1.001  4000
fitnew   20.043  3.422  14.100  17.630  19.770  22.292  27.360 1.001  4000
deviance 446.255 12.290 424.000 437.700 445.600 454.100 472.302 1.001  4000

For each parameter, n.eff is a crude measure of effective sample size,
and Rhat is the potential scale reduction factor (at convergence, Rhat=1).

DIC info (using the rule, pD = Dbar-Dhat)
pD = 66.0 and DIC = 512.2
DIC is an estimate of expected predictive error (lower deviance is better).
>

```

The Bayesian p-value for this model is

```

> mean(out$sims.list$fit>out$sims.list$fitnew)
[1] 0.473
>

```

indicating a pretty good fit. Given the site-level random effect, it would be surprising for this model to not fit! One thing we notice is that the posterior standard deviations of the regression parameters are much higher, a result of the excess variation. (we would also notice much less precise predictions of hypothetical new observations).

1.11 Binomial GLMs

Another class of statistical models that are very important in ecology are binomial models. We use binomial models for count data whenever the observations are counts or frequencies and it is natural to condition on a “sample size” - the maximum frequency possible in a sample, say K (i.e., K is known). The random variable, y/leK , is then the frequency of occurrences out of K . The parameter of the binomial models is p , often called “success probability” which is related to the expected value of y by $E[y] = pK$. Binomial GLMs or binomial regression models are often referred to as logistic regression, but that term really

only applies when the logistic link is used to model the relationship between p and covariates (see below).

One of the most typical Binomial GLMs occurs when the sample size equals 1 and the outcome, y , is “presence” ($y = 1$) or “absence” ($y = 0$) of a species. This is a classical “species distribution” modeling situation. A special situation occurs when presence/absence is observed with error (MacKenzie et al., 2002; MacKenzie, 2006; Kéry et al., 2010). In that case, $K > 1$ samples are usually required in order to estimate model parameters effectively. In standard binomial regression problems the sample size is fixed by design but interesting models also arise when the sample size is itself a random variable. These are the N-mixture models (Royle, 2004; Kéry et al., 2005; Royle and Dorazio, 2008; Kéry, 2010) ch. 22) and related models (in this case, N being the sample size which we labeled K above). This is actually a little bit confusing because the binomial index is usually referred to as “sample size” but in this context N is actually a “population size”. A useful situation in which the binomial sample size is “fixed” is closed population capture-recapture models in which a population of individuals is sampled K times. The number of times each individual is encountered is a binomial outcome with parameter - encounter probability - p , based on a sample of size K . We consider such models in the following chapter.

1.11.1 Binomial regression

In binomial models, covariates are modeled on a suitable transformation (the link function) of the binomial success probability, p . Let x_i denote some measured covariate for sample unit i and let p_i be the success probability for unit i . The standard choice is the “logit” link function which is:

$$\log(p[i]/(1 - p[i])) = \alpha + \beta * x[i]$$

with inverse “expit”

$$p[i] = \text{expit}(\alpha + \beta * x[i]) = \exp(\alpha + \beta * x[i]) / (1 + \exp(\alpha + \beta * x[i]))$$

There are many other possible link functions. However, ecologists seem to blindly adopt the logit link function without question to such an extent that you are likely to be questioned by referees and associate editors if you use some alternative link (unless you are doing species distribution modeling, in which case any explicit link function will be questioned by some referees). We sometimes use the “complementary log-log” (= “cloglog”) link function in ecological applications because it can often be justified based on subject-matter considerations (Royle and Dorazio (2008); section XYZ) or natural scaling relationships germane to the problem. For example, the cloglog link arises as the “probability of a count greater than 0” under a Poisson model. That is, $\Pr(y > 0) = 1 - \exp(-\lambda)$ in which case

$$\text{cloglog}(p) = \log(-\log(1 - p)) = \log(\lambda)$$

So that if you have covariates in your linear predictor for $E[y]$ under a Poisson model then they are linear on the complementary log-log link of p . We will use the cloglog link in some analyses of SCR models in Chapter 4 and elsewhere.

A natural situation in which the cloglog link arises is modeling occupancy in which $N \sim \text{Poisson}(A * \lambda)$ and you have site area, A , measured for every sample. In this case the probability that the site is occupied, ψ , is related to area on the cloglog scale. i.e.,

$$\text{cloglog}(\psi) = \log(A) + \log(\lambda).$$

There seems to be perennial debate over whether site area should be a covariate on “detection” or “occupancy” and the above argument suggests the latter.

1.11.2 Example: Waterfowl Banding Data

It would be easy to consider a standard “distribution modeling” application where $K = 1$ and the outcome is occurrence ($y = 1$) or not ($y = 0$) of some species. Such examples abound in books (e.g., Royle and Dorazio (2008), ch. 3; Kéry (2010), chapter 21 XYZ?; Kéry and Schaub (2011), chapter XYZ) and in the literature (see Kéry et al. (2010); Kéry et al. (2010) XYZ). Instead, we will consider an example involving band returns of waterfowl which were analyzed by Royle and Dubovsky (200X)⁶.

For these data, $y[i]$ is the number of waterfowl bands recovered out of $B[i]$ birds banded at some location $s[i]$. In this case $B[i]$ is fixed. Thinking about recovery rate as being proportional to harvest rate, we wanted to explore geographic gradients in recovery rate resulting from variability in harvest pressure experienced by populations depending on their migration ecology. As such, we fit a basic binomial GLM with a linear response to geographic coordinates (including an interaction term). The data are provided on the web supplement along with an R script to do the post-processing. Here we just provide the part of the script for creating the model and calling WinBUGS:

```

sink("model.txt")
cat("
model {
  for(t in 1:5){
    for (i in 1:nobs){
      m[i,t] ~ dbin(p[i,t], R[i,t])
      logit(p[i,t]) <- alpha0[t] + alpha1*X[i,1] + alpha2*X[i,2] + alpha3*X[i,1]*X[i,2]
    }
  }
  alpha1~dnorm(0,.001)
  alpha2~dnorm(0,.001)
  alpha3~dnorm(0,.001)
  for(t in 1:5){
    alpha0[t] ~ dnorm(0,.001)
  }
}
```

⁶not happy about this example. Anyone got a better one?

```

1213 ",fill=TRUE)
1214 sink()
1215
1216 data <- list('R', 'm', 'nobs','X')
1217 inits <- function(){
1218   list(alpha0=rnorm(5),alpha1=0,alpha2=0,alpha3=0)
1219 }
1220 parms <- list('alpha0','alpha1','alpha2','alpha3')
1221 out <- bugs(data,inits, parms,"model.txt",n.chains=3,
1222   n.iter=2000,n.burnin=1000,
1223   n.thin=2, debug=TRUE)

```

1224 Posterior summaries of model parameters are as follows:

```

1225 Inference for Bugs model at "model.txt", fit using WinBUGS,
1226   3 chains, each with 2000 iterations (first 1000 discarded), n.thin = 2
1227   n.sims = 1500 iterations saved
1228
1229      mean      sd    2.5%    25%    50%    75%    97.5%  Rhat  n.eff
1230 alpha0[1] -2.346 0.036  -2.417  -2.370  -2.346  -2.323  -2.277 1.001 1500
1231 alpha0[2] -2.356 0.032  -2.420  -2.379  -2.356  -2.335  -2.292 1.001 1500
1232 alpha0[3] -2.220 0.035  -2.291  -2.244  -2.219  -2.197  -2.153 1.001 1500
1233 alpha0[4] -2.144 0.039  -2.225  -2.169  -2.143  -2.116  -2.068 1.000 1500
1234 alpha0[5] -1.925 0.034  -1.990  -1.949  -1.924  -1.901  -1.856 1.004   570
1235 alpha1     -0.023 0.003   -0.028   -0.025   -0.023   -0.022   -0.018 1.001 1500
1236 alpha2      0.020 0.006    0.009    0.016    0.020    0.024    0.031 1.001 1500
1237 alpha3      0.000 0.001   -0.002   -0.001    0.000    0.000    0.002 1.001 1500
1238 deviance 1716.001 4.091 1710.000 1713.000 1715.000 1718.000 1726.000 1.001 1500

```

1239 For each parameter, n.eff is a crude measure of effective sample size,
 1240 and Rhat is the potential scale reduction factor (at convergence, Rhat=1).

1241
 1242 DIC info (using the rule, $pD = \bar{D} - \hat{D}$)

1243 $pD = 7.9$ and $DIC = 1723.9$

1244 DIC is an estimate of expected predictive error (lower deviance is better).

1245 The basic result suggests a negative east-west gradient and a positive south
 1246 to north gradient but no interaction. A map of the response surface is given
 1247 below. We could use DIC to do some model selection - i.e., try models with-
 1248 out the interaction term, or models with a quadratic term, or with a constant
 1249 intercept, etc., but we don't pursue that here. We did an MCMC run where
 1250 we saved the binomial parameter p and computed the Bayesian p -value [double
 1251 use of " p " here is confusing!] using a fit statistic based on the Freeman-Tukey
 1252 statistic (see Section XXX above). The result indicates that the linear response
 1253 surface model does not provide an adequate fit of the data. The reader should
 1254 contemplate whether this invalidates the basic interpretation of the result.

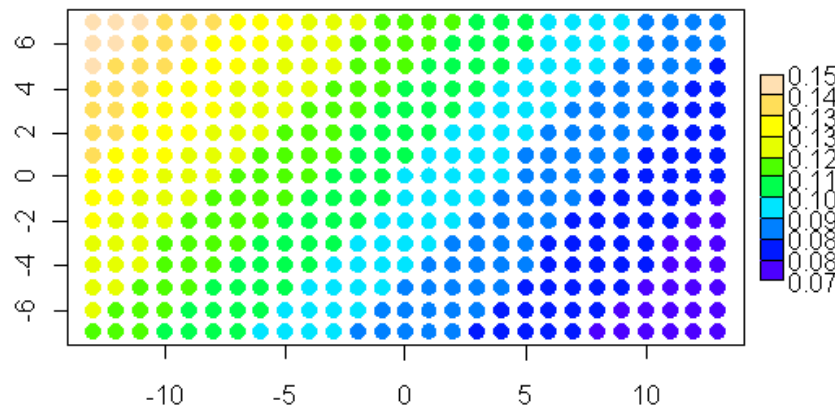


Figure 1.6: Needs a caption

1.12 Stuff on hierarchical models here?

1.13 Summary and Outlook

GLMs and GLMMs are the most useful statistical methods in all of ecology. The principles and procedures underlying these methods are relevant to nearly all modeling and analysis problems in every branch of ecology. Moreover, understanding how to analyze these models is crucial in a huge number of diverse problems. If you understand and can conduct classical likelihood and Bayesian analysis of Poisson and binomial GLM(M)s, then you will be successful analyzing and understanding more complex classes of models that arise. We will see shortly that spatial capture-recapture models are just a type of GLMM (i.e., a GLM with a random effect) and thus having a basic understanding of the conceptual origins and formulation of GLMs and their analysis is extremely useful. We note that GLMs are routinely analyzed by likelihood methods but we have focused on Bayesian analysis here in order to develop the tools that are less familiar to most ecologists. In particular, Bayesian analysis of GLMs with random effects (i.e., GLMMs) is relatively straightforward because the models are easy to analyze conditional on the random effect, using methods of MCMC. Thus, we will often analyze SCR models in later chapters by MCMC, explicitly adopting a Bayesian inference framework.

In that regard, BUGS engines are enormously useful because they provides a straightforward way to carry out analyses by MCMC by just describing the model, and not having to worry about how to actually build MCMC algorithms.

1277 That said, the BUGS language is more important than just to the extent that
1278 it enables one to do MCMC - it is useful as a modeling tool because it fosters
1279 understanding, in the sense that it forces you to become intimate with your
1280 model. You have to write down all of the probability assumptions, the relation-
1281 ships between observations and latent variables and parameters. This is really
1282 a great learning paradigm that you can grow with. Skills gained in Bayesian
1283 analysis of the GLMMs covered in this chapter will be directly transferrable and
1284 useful for the SCR models addressed subsequently. Before getting to that, how-
1285 ever, it will be useful to talk about more basic, conventional closed population
1286 capture-recapture models and these are the topic of the next Chapter.

Bibliography

- Brooks, S. P., Catchpole, E. A., and Morgan, B. J. T. (2000), “Bayesian Animal Survival Estimation,” *Statistical Science*, 15, 357–376.
- Chandler, R. and Royle, J. (2012), “Spatially-explicit models for inference about density in unmarked populations,” *Biometrics (in review)*.
- Gardner, B., Royle, J., Wegan, M., Rainbolt, R., and Curtis, P. (2010), “Estimating black bear density using DNA data from hair snares,” *The Journal of Wildlife Management*, 74, 318–325.
- Gelman, A., Carlin, J. B., Stern, H. S., and Rubin, D. B. (2004), *Bayesian data analysis, second edition.*, Boca Raton, Florida, USA: CRC/Chapman & Hall.
- Gelman, A., Meng, X. L., and Stern, H. (1996), “Posterior predictive assessment of model fitness via realized discrepancies,” *Statistica Sinica*, 6, 733–759.
- Hawkins, C. and Racey, P. (2005), “Low population density of a tropical forest carnivore, *Cryptoprocta ferox*: implications for protected area management,” *Oryx*, 39, 35–43.
- Jackson, R., Roe, J., Wangchuk, R., and Hunter, D. (2006), “Estimating Snow Leopard Population Abundance Using Photography and Capture-Recapture Techniques,” *Wildlife Society Bulletin*, 34, 772–781.
- Kéry, M. (2010), *Introduction to WinBUGS for Ecologists: Bayesian Approach to Regression, ANOVA, Mixed Models and Related Analyses*, Academic Press.
- Kéry, M., Gardner, B., Stoeckle, T., Weber, D., and Royle, J. A. (2010), “Use of Spatial Capture-Recapture Modeling and DNA Data to Estimate Densities of Elusive Animals,” *Conservation Biology*, 25, 356–364.
- Kéry, M., Royle, J., and Schmid, H. (2005), “Modeling avian abundance from replicated counts using binomial mixture models,” *Ecological Applications*, 15, 1450–1461.
- Kery, M. and Schaub, M. (2011), *Bayesian Population Analysis Using WinBugs*, Academic Press.

- 1316 King, R. (2009), “Missing,” *missing*, Missing.
- 1317 Le Cam, L. (1990), “Maximum likelihood: an introduction,” *International Sta-*
1318 *tistical Review/Revue Internationale de Statistique*, 153–171.
- 1319 Link, W. A. and Barker, R. J. (2009), *Bayesian Inference: With Ecological*
1320 *Applications*, London, UK: Academic Press.
- 1321 MacEachern, S. and Berliner, L. (1994), “Subsampling the Gibbs sampler,”
1322 *American Statistician*, 188–190.
- 1323 MacKenzie, D. (2006), *Occupancy estimation and modeling: inferring patterns*
1324 *and dynamics of species occurrence*, Academic Press.
- 1325 MacKenzie, D. I., Nichols, J. D., Lachman, G. B., Droege, S., Royle, J. A., and
1326 Langtimm, C. A. (2002), “Estimating site occupancy rates when detection
1327 probabilities are less than one,” *Ecology*, 83, 2248–2255.
- 1328 McCarthy, M. A. (2007), *Bayesian Methods for Ecology*, Cambridge: Cambridge
1329 University Press.
- 1330 McCullagh, P. and Nelder, J. (1989), *Generalized linear models*, Chapman &
1331 Hall/CRC.
- 1332 Nelder, J. and Wedderburn, R. (1972), “Generalized linear models,” *Journal of*
1333 *the Royal Statistical Society. Series A (General)*, 370–384.
- 1334 Royle, J. (2008), “Analysis of capture-recapture models with individual covari-
1335 ates using data augmentation,” *Biometrics*, 267–274.
- 1336 Royle, J. and Dorazio, R. (2008), *Hierarchical modeling and inference in ecol-*
1337 *ogy: the analysis of data from populations, metapopulations and communities*,
1338 Academic Press.
- 1339 Royle, J. and Link, W. (2006), “Generalized site occupancy models allowing for
1340 false positive and false negative errors,” *Ecology*, 87, 835–841.
- 1341 Royle, J. A. (2004), “Generalized estimators of avian abundance from count
1342 survey data,” *Animal Biodiversity and Conservation*, 27, 375–386.
- 1343 Sepúlveda, M., Bartheld, J., Monsalve, R., Gómez, V., and Medina-Vogel, G.
1344 (2007), “Habitat use and spatial behaviour of the endangered Southern river
1345 otter (*Lontra provocax*) in riparian habitats of Chile: conservation implica-
1346 tions,” *Biological Conservation*, 140, 329–338.
- 1347 Sturtz, S., Ligges, U., and Gelman, A. (2005), “R2WinBUGS: A Package for
1348 Running WinBUGS from R,” *Journal of Statistical Software*, 12, 1–16.
- 1349 Trolle, M. and Kéry, M. (2005), “Camera-trap study of ocelot and other secretive
1350 mammals in the northern Pantanal,” *Mammalia*, 69, 409–416.
- 1351 Zuur, A., Ieno, E., Walker, N., Saveliev, A., and Smith, G. (2009), *Mixed effects*
1352 *models and extensions in ecology with R*, Springer Verlag.