

Chapter 1

Developing Markov Chain Monte Carlo Samplers

In this chapter we will dive a little deeper into Markov chain Monte Carlo (MCMC) sampling. We will construct custom MCMC samplers in **R**, starting with easy-to-code GLMs and GLMMs and moving on to simple CR and SCR models. This material might seem slightly out of place here, as it does not deal with specific aspects or modifications of SCR models, but rather, with a particular way of implementing them (and other models, too). Knowing how to build an MCMC sampler is not essential for any of the SCR models we have covered so far, but we will need these skills to implement some models that come up in the last few chapters of this book. The aim of this chapter is to provide you with some working knowledge of building MCMC samplers. To this end, we will NOT provide exhaustive background information on the theory and justification of MCMC sampling – there are entire books dedicated to that subject and we refer you to Robert and Casella (2004) and Robert and Casella (2010). Rather we aim to provide you with enough background and technical know-how to start building your own MCMC samplers for SCR models in **R**. You will find that quite a few topics that come up in this chapter have already been covered in previous chapters, particularly the introduction into Bayesian analysis in Chapt. ???. To keep you from having to leaf back and forth we will in some places briefly review aspects of Bayesian analysis, but we try to focus on the more technical issues of building MCMC samplers relevant to SCR models.

1.1 Why Build Your Own MCMC Algorithm?

The standard programs we have used so far to do MCMC analyses are **WinBUGS** (Gilks et al., 1994) and **JAGS** (Plummer, 2003). The wonderful thing about these **BUGS** engines is that they automatically use appropriate and,

most of the time, reasonably efficient forms of MCMC sampling for the model specified by the user.

The fact that we have such a Swiss Army knife type of MCMC machine begs the question: Why would anyone want to build their own MCMC algorithm? For one, there are a limited number of distributions and functions implemented in **BUGS**. While **OpenBUGS** and **JAGS** provide more options, some more complex models may be impossible to build within these programs. A very simple example from spatial capture-recapture that can give you a headache in **WinBUGS** is when your state-space is an irregular-shaped polygon, rather than an ideal rectangle that can be characterized by four pairs of coordinates. It is easy to restrict activity centers to any arbitrary polygon in **R** using an ESRI shapefile (and we will show you an example in a little bit), but you cannot use a shapefile in a **BUGS** model. Similarly, models of space usage that take into account ecological distance (Chapt. ??) cannot be implemented in the **BUGS** engines.

Sometimes, implementing an MCMC algorithm in **R** may be faster than in **WinBUGS** - especially if you want to run simulation studies where you have hundreds or more simulated data sets, several years' worth of data or other large models, this can be a big advantage. Further, writing your own sampler gives you more control over which kind of updater is used (see following sections). Finally, building your own MCMC algorithm is a great exercise to understand how MCMC sampling works. So while using the **BUGS** language requires you to understand the structure of your model, building an MCMC algorithm requires you to think about the relationship between your data, priors and posteriors, and how these can be efficiently analyzed and characterized. However, if you don't think you will ever sit down and write your own MCMC sampler, consider skipping this chapter - apart from coding it will not cover anything SCR-related that is not covered by other, more model-oriented chapters as well.

1.2 MCMC and Posterior Distributions

MCMC is a class of simulation methods for drawing (correlated) random numbers from a target distribution, which in Bayesian inference is the posterior distribution. As a reminder, the posterior distribution is a probability distribution for an unknown parameter, say θ , given observed data and its prior probability distribution (the probability distribution we assign to a parameter before we observe data). The great benefit of having the posterior distribution of θ is that it can be used to make probability statements about θ , such as the probability that θ is equal to some value, or the probability that θ falls within some range of values. The posterior distribution summarizes all we know about a parameter and thus, is the central object of interest in Bayesian analysis. Unfortunately, in many if not most practical applications, it is nearly impossible to directly compute the posterior. Recall Bayes' theorem:

$$[\theta|y] = \frac{[y|\theta][\theta]}{[y]}, \quad (1.1)$$

where θ is the parameter of interest, y is the observed data, $[\theta|y]$ is the posterior, $[y|\theta]$ the likelihood of the data conditional on θ , $[\theta]$ the prior probability of θ , and, finally, $[y]$ is the marginal probability of the data, defined as

$$[y] = \int [y|\theta][\theta]d\theta$$

This marginal probability is a normalizing constant that ensures that the posterior integrates to 1. Often, the integral is difficult or impossible to evaluate, unless you are dealing with a really simple model. For example, consider a normal model, with a set of n observations, $y_i; i = 1, 2, \dots, n$:

$$y_i \sim \text{Normal}(\mu, \sigma),$$

where σ is known and our objective is to estimate μ . To fully specify the model in a Bayesian framework, we first have to define a prior distribution for μ . Recall from Chapt. ?? that for certain data models, certain priors lead to conjugacy, i.e. if you choose a certain prior for your parameter, the posterior distribution will be of a known parametric form. More specifically, under conjugacy, the prior and posterior distributions are from the same parametric family. The conjugate prior for the mean of a normal model is also a normal distribution:

$$\mu \sim \text{Normal}(\mu_0, \sigma_0^2).$$

If μ_0 and σ_0^2 are fixed, the posterior for μ has the following form (for some of the algebra behind this, see Chapt. 2 in Gelman et al. (2004)):

$$\mu|y \sim \text{Normal}(\mu_n, \sigma_n^2) \tag{1.2}$$

where

$$\mu_n = \left(\frac{\sigma^2}{\sigma^2 + n\sigma_0^2} \right) \times \left(\mu_0 + \frac{n\sigma_0^2}{\sigma^2 + n\sigma_0^2} \right) \times \bar{y}$$

and

$$\sigma_n^2 = \frac{\sigma^2 \sigma_0^2}{\sigma^2 + n\sigma_0^2}.$$

We can directly obtain estimates of interest from this normal posterior distribution, such as its mean $\hat{\mu}$ (which is equivalent to an estimate of μ_n) and variance; we do not need to apply MCMC, since we can recognize the posterior as a parametric distribution, including the normalizing constant $[y]$. But generally we will be interested in more complex models with several, say m , parameters. In this case, computing $[y]$ from Eq. 1.1 requires m -dimensional integration, which can be difficult or impossible. Thus, the posterior distribution is generally only known up to a constant of proportionality:

$$[\theta|y] \propto [y|\theta][\theta]$$

The power of MCMC is that it allows us to approximate the posterior using simulation without evaluating the high dimensional integrals, and to directly sample

from the posterior, even when the posterior distribution is unknown! The price is that MCMC is computationally expensive. Although MCMC first appeared in the scientific literature in 1949 (Metropolis and Ulam, 1949), widespread use did not occur until the 1980s when computational power and speed increased (Gelfand and Smith, 1990). It is safe to say that the advent of practical MCMC methods is the primary reason why Bayesian inference has become so popular during the past three decades.

In a nutshell, MCMC lets us generate sequential draws of θ (the parameter(s) of interest) from distributions approximating the unknown posterior over T iterations. The distribution of the draw at t depends on the value drawn at $t-1$; hence, the draws from a Markov chain¹. As T goes to infinity, the Markov chain converges to the desired distribution, in our case the posterior distribution for $\theta|y$. Thus, once the Markov chain has reached its stationary distribution, the generated samples can be used to characterize the posterior distribution, $[\theta|y]$, and point estimates of θ , its standard error and confidence bounds, can be obtained directly from this approximation of the posterior.

1.3 Types of MCMC Sampling

There are several general MCMC algorithms in widespread use, the most popular being Gibbs sampling and Metropolis-Hastings sampling, both of which were briefly introduced in Chapt. ???. We will be dealing with these two classes in more detail and use them to construct MCMC algorithms for SCR models. Also, we will briefly review alternative techniques that are applicable in some situations.

1.3.1 Gibbs sampling

Gibbs sampling was named after the physicist J.W. Gibbs by Geman and Geman (1984), who applied the algorithm to a Gibbs distribution². The roots of Gibbs sampling can be traced back to work of Metropolis et al. (1953), and it is actually closely related to Metropolis sampling (see Chapt. 11.5 in Gelman et al. (2004), for the link between the two samplers). We will focus on the technical aspects of this algorithm, but if you find yourself hungry for more background, Casella and George (1992) provide a more in-depth introduction to the Gibbs sampler.

Let's go back to our simple example from above to understand the motivation and functioning of Gibbs sampling. Recall that for a normal model with known variance and a normal prior for μ , the posterior distribution of $\mu|y$ is also normal. Conversely, with a fixed (known) μ , but unknown variance, the conjugate prior for σ^2 is an inverse-gamma distribution with shape and scale parameters a and

¹Remember that for T random samples $\theta^{(1)}, \dots, \theta^{(T)}$ from a Markov chain the distribution of $\theta^{(t)}$ depends only on the immediately preceding value, $\theta^{(t-1)}$.

²a distribution from physics we are not going to worry about, since it has no immediate connection with Gibbs sampling other than giving its name

133 b :

$$\sigma^2 \sim \text{Inverse-Gamma}(a, b).$$

134 With fixed a and b , algebra reveals that the posterior $[\sigma^2|\mu, y]$ is also an inverse-
135 gamma distribution, namely:

$$\sigma^2|\mu, y \sim \text{Inverse-Gamma}(a_n, b_n), \quad (1.3)$$

136 where $a_n = n/2 + a$ and $b_n = (1/2) \sum_{i=1}^n (y_i - \mu)^2 + b$. However, what if we know
137 neither μ nor σ^2 , which is probably the more common case? The joint posterior
138 distribution of μ and σ^2 now has the general structure

$$[\mu, \sigma^2|y] = \frac{[y|\mu, \sigma^2][\mu][\sigma^2]}{\int [y|\mu][\mu][\sigma^2] d\mu d\sigma^2}$$

139 or

$$[\mu, \sigma^2|y] \propto [y|\mu, \sigma^2][\mu][\sigma^2]$$

140 This cannot easily be reduced to a distribution we recognize. However, we
141 can condition μ on σ^2 (i.e., we treat σ^2 as fixed) and remove all terms from the
142 joint posterior distribution that do not involve μ to construct the full conditional
143 distribution,

$$[\mu|\sigma^2, y] \propto [y|\mu][\mu]$$

144 The full conditional of μ again takes the form of the normal distribution
145 shown in Eq. 1.2; similarly, $[\sigma^2|\mu, y]$ takes the form of the inverse-gamma dis-
146 tribution shown in Eq. 1.3, both distributions we can easily sample from. And
147 this is precisely what we do when using Gibbs sampling: we break down high-
148 dimensional problems into convenient one-dimensional problems by constructing
149 the full conditional distributions for each model parameter separately; and we
150 sample from these full conditionals, which, if we choose conjugate priors, are
151 known parametric distributions. Let's put the concept of Gibbs sampling into
152 the MCMC framework of generating successive samples, using our simple nor-
153 mal model with unknown μ and σ^2 and conjugate priors as an example. These
154 are the steps you need in order to build a Gibbs sampler:

155 **Step 0:** Begin with some initial values for θ , say $\theta^{(0)}$. In our example, $\theta = (\mu, \sigma)$,
156 so we have to specify initial values for μ and σ , for example by drawing a random
157 number from some uniform distribution, or by setting them close to what we
158 think they might be. (Note: This step is required in any MCMC sampling;
159 chains have to start from somewhere. We will get back to these technical details
160 a little later.)

161 **Step 1:** For iteration t , Draw $\theta^{(t)}$ from the conditional distribution $[\theta_1^{(t)}|\theta_2^{(t-1)}, \dots,$
162 $\theta_d^{(t-1)}]$. Here, θ_1 is μ , which we draw from the normal distribution in Eq. 1.2
163 using $\sigma^{(t-1)}$ as value for σ .

164 **Step 2:** Draw $\theta_2^{(t)}$ from the conditional distribution $[\theta_2^{(t)}|\theta_1^{(t)}, \theta_3^{(t-1)}, \dots, \theta_d^{(t-1)}]$.
 165 Here, θ_2 is σ , which we draw from the inverse-gamma distribution of Eq. 1.3,
 166 using the newly generated $\mu^{(t)}$ as value for μ .

167 **Step 3, ..., d:** Draw $\theta_3^{(t)}, \theta_4^{(t)}, \dots, \theta_d^{(t)}$ from their conditional distribution
 168 $[\theta_3^{(t)}|\theta_1^{(t)}, \theta_2^{(t)}, \theta_4^{(t-1)}, \dots, \theta_d^{(t-1)}], \dots, [\theta_d^{(t)}|\theta_1^{(t)}, \dots, \theta_{d-1}^{(t)}]$. In our example we
 169 have no additional parameters, so we only need step 0 through to 2.

170 **Repeat Steps 1 to d** for $T =$ a large number of samples.

171 Note that the order in which we update the parameters within the Gibbs
 172 algorithm does not matter. In terms of **R** coding, this means we have to write
 173 Gibbs updaters for μ and σ^2 and embed them into a loop over T iterations. The
 final code in the form of an **R** function is shown in Panel 1.1.

```

Norm.Gibbs<-function(y=y,mu_0=mu_0,sigma2_0=sigma2_0,a=a,b=b,niter=niter){

  ybar<-mean(y)
  n<-length(y)
  mu<-1          #mean initial value
  sigma2<-1      #sigma2 initial value
  an<-n/2 + a    #shape parameter of IvGamma of sigma2
  out<-matrix(nrow=niter, ncol=2)
  colnames(out)<-c('mu', 'sig')

  for (i in 1:niter) {

    #update mu
    mu_n<-((sigma2/(sigma2+n*sigma2_0))*mu_0
    + (n*sigma2_0/(sigma2 + n*sigma2_0))*ybar)
    sigma2_n <- (sigma2*sigma2_0)/ (sigma2 + n*sigma2_0)
    mu<-rnorm(1,mu_n, sqrt(sigma2_n))

    #update sigma2
    bn<- 0.5 * (sum((y-mu)^2)) + b
    sigma2<-1/rgamma(1,shape=an, rate=bn)
    out[i,<-c(mu,sqrt(sigma2))
  }
  return(out)
}

```

Panel 1.1: R-code for a Gibbs sampler for a normal model with unknown μ and σ and conjugate priors (normal and inverse-gamma, respectively) for both parameters.

175 This is it! You can go ahead and simulate some data, $y \sim \text{Normal}(5, 0.5)$
 176 and then use the function `NormGibbs()` in the **R** package `scrbook` to run your
 177 first Gibbs sampler (note that the **R** function `rnorm` requires you to supply the
 178 standard deviation σ and we have written `NormGibbs` so that it returns σ instead
 179 of σ^2 so you can easily compare you input value and parameter estimate).

```
180 > set.seed(13)
181
182 #true mean and sd are 5 and 0.5
183 > y<-rnorm(1000, 5,0.5) #data
184
185 > mu_0<-0 #prior mean
186 > sigma2_0<-100 #prior variance
187
188 #inverse-gamma hyperparameters
189 > a<-0.1
190 > b<-0.1
191
192 > mod=Norm.Gibbs(y, mu_0, sigma2_0, a,b,niter=10000)
```

193 Your output, `mod`, will be a table with two columns, one per parameter, and
 194 T rows, one per iteration. For this 2-parameter example you can visualize the
 195 joint posterior by plotting samples of μ against samples of σ (Fig. 1.1):

```
196 > plot(out[,1], out[,2])
```

197 The marginal distribution of each parameter is approximated by examining
 198 the samples of this particular parameter. You can visualize it by plotting a
 199 histogram of the samples (Fig. 1.2 upper left and right):

```
200 > par(mfrow=c(1,2))
201 > hist(out[,1]); hist(out[,2])
```

202 Finally, recall an important characteristic of MCMC, namely, that the chain
 203 has to have converged (reached its stationary distribution) in order to regard
 204 samples as coming from the posterior distribution. In practice, that means you
 205 have to throw out some of the initial samples called the burn-in. We will talk
 206 about this in more detail when we talk about convergence diagnostics. For now,
 207 you can use the `plot(out[,1])` or `plot(out[,2])` command to make a time
 208 series plot of the samples of each parameter and visually assess how many of the
 209 initial samples you should discard. Fig. 1.2 bottom left and right shows plots
 210 for the samples of μ and σ from our simulated data set; you see that in this
 211 simple example the Markov chain apparently reaches its stationary distribution
 212 very quickly – the chains look ‘grassy’ seemingly from the start. It is hard to
 213 discern a burn-in phase visually (but we will see examples further on where the
 214 burn-in is clearer) and you may just discard the first 500 draws to be sure you
 215 only use samples from the posterior distribution. The mean of the remaining
 216 samples are your estimates of μ and σ :

```

217 > summary(mod[501:10000,])
218      mu      sig
219 Min.   :4.935  Min.   :0.4652
220 1st Qu.:4.988  1st Qu.:0.4930
221 Median :4.998  Median :0.5006
222 Mean   :4.998  Mean    :0.5008
223 3rd Qu.:5.009  3rd Qu.:0.5084
224 Max.   :5.062  Max.    :0.5486

```

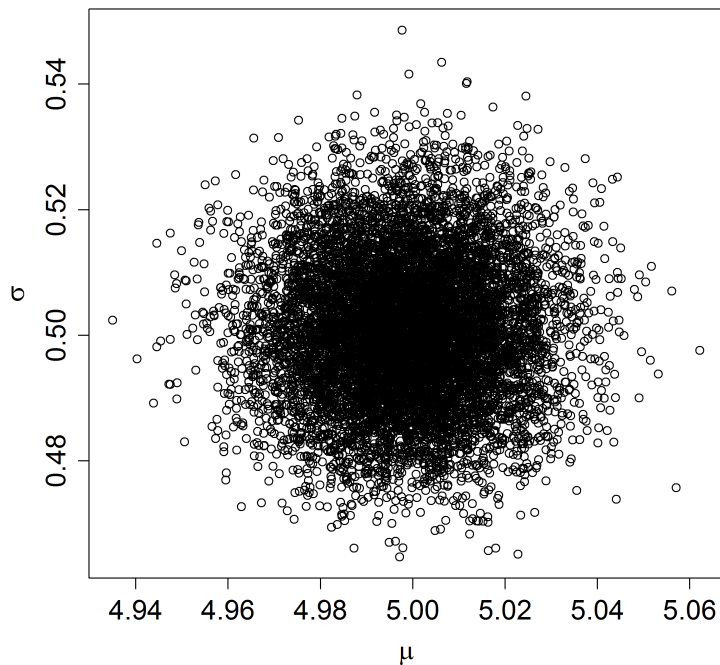


Figure 1.1: Joint posterior distribution of μ and σ from a normal Model

225 1.3.2 Metropolis-Hastings sampling

226 Although it is applicable to a wide range of problems, the limitations of Gibbs
 227 sampling are obvious: what if we do not want to use conjugate priors or what
 228 if we cannot recognize the full conditional distribution as a parametric distri-
 229 bution, or simply do not want to worry about these issues? The most general
 230 solution is to use the Metropolis-Hastings (MH) algorithm, which also goes back

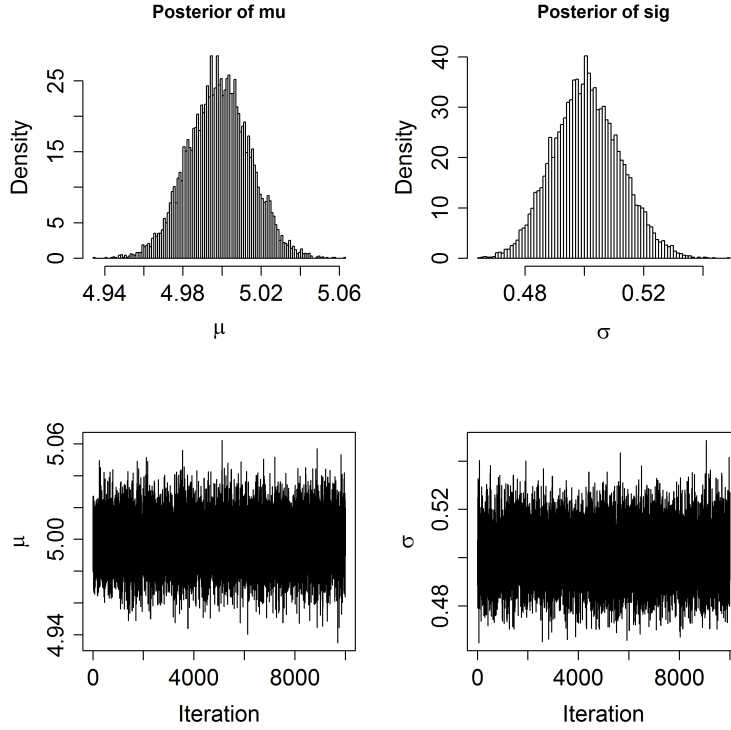


Figure 1.2: Plots of the posterior distributions of μ (upper left) and σ (upper right) from a normal model and time series plots of μ (lower left) and σ (lower right).

231 to the work by Metropolis et al. (1953). You saw the basics of this algorithm
 232 in Chapt. ???. In a nutshell, because we do not recognize the posterior $[\theta|y]$ as
 233 a parametric distribution, the MH algorithm generates samples from a known
 234 proposal distribution, say $h(\theta)$, that depends on the value of θ at the previous
 235 time step, $\theta^{(t-1)}$. The candidate value θ^* is accepted with probability

$$r = \min\left(1, \frac{[\theta^*|y]h(\theta^{(t-1)}|\theta^*)}{[\theta^{(t-1)}|y]h(\theta^*|\theta^{(t-1)})}\right)$$

236 Proposal distributions must be chosen so that reversibility is ensured. That
 237 means, it must be possible to go from any one value to any other. But within
 238 that criterion the proposal distribution can be absolutely anything! You can
 239 generate candidate values from a Normal(0,1) distribution, from a Uniform(-
 240 3455,3455) distribution, or anything of proper support. Note, however, that
 241 good choices of $h(\theta)$ are those that approximate the posterior distribution. Ob-
 242 viously if $h(\theta) = [\theta|y]$ (i.e., the posterior) then you always accept the draw

and it stands to reason that proposals that are more similar to $[\theta|y]$ will lead to higher acceptance probabilities. Actually, when $h(\theta) = [\theta|y]$ we can draw samples of θ directly from $h(\theta)$, which brings us back to Gibbs sampling. Thus, Gibbs sampling is a special case of Metropolis-Hastings sampling.

The original Metropolis algorithm required $h(\theta)$ to be symmetric so that

$$h(\theta^*|\theta^{(t-1)}) = h(\theta^{(t-1)}|\theta^*)$$

In that case these two terms just cancel out from the MH acceptance probability and r is then just the ratio of the target density evaluated at the candidate value to that evaluated at the current value. A later development of the algorithm by Hastings (1970) lifted this condition. Since using a symmetric proposal distribution makes life a little easier, we are going to focus on this specific case. 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 δ , which is prescribed by the user (see below for further explanation).

Parameters with bounded support: Many models contain parameters that have bounded support. E.g., variance parameters live on $[0, \infty]$, parameters that represent probabilities live on $[0, 1]$, etc.. For such cases, it is sometimes convenient to use a random walk proposal distribution that can generate any real number (e.g., a normal random walk proposal). Under these circumstances you should not constrain the proposal distribution itself, but you can just reject parameters that are outside of the parameter space (sec. 6.4.1 in Robert and Casella, 2010). You will see plenty of examples of updating parameters with bounded support in this chapter.

It is worth knowing that there are alternatives to the random walk MH algorithm. For example, in the independent MH, the proposal distribution h does not depend on $\theta^{(t-1)}$, while the Langevin algorithm (Roberts and Rosenthal, 1998) aims at avoiding the random walk by favoring moves towards regions of higher posterior probability density. The interested reader should look up these algorithms in Robert and Casella (2004) or Robert and Casella (2010).

Building a MH sampler can be broken down into several steps. We are going to demonstrate these steps using a different but still simple and common model: the logit-normal or logistic regression model. For simplicity, assume that

$$y|\theta \sim \text{Bernoulli}\left(\frac{\exp(\theta)}{1 + \exp(\theta)}\right)$$

and

$$\theta \sim \text{Normal}(\mu, \sigma).$$

The following steps are required to set up a random walk MH algorithm:

Step 0: Choose initial values, $\theta^{(0)}$.

Step 1: Generate a proposed value of θ from $h(\theta^*|\theta^{(t-1)})$. We will use the random walk MH algorithm, so we draw θ^* from $\text{Normal}(\theta^{(t-1)}, \delta)$, where δ is the

standard deviation of the normal proposal distribution, the tuning parameter that we have to set.

Step 2: Calculate the ratio of posterior densities for the proposed and the original value for θ :

$$r = \frac{[\theta^*|y]}{[\theta^{(t-1)}|y]}.$$

In our example,

$$r = \frac{\text{Bernoulli}(y|\theta^*) \times \text{Normal}(\theta^*|\mu, \sigma)}{\text{Bernoulli}(y|\theta^{(t-1)}) \times \text{Normal}(\theta^{(t-1)}|\mu, \sigma)}$$

Step 3: Set

$$\begin{aligned}\theta^t &= \theta^* \text{ with probability } \min(r, 1) \\ &= \theta^{(t-1)} \text{ otherwise}\end{aligned}$$

We can do this last step by drawing a random number u from a $\text{Uniform}(0, 1)$ and accept θ^* if $u < r$. This is repeated for $t = 1, 2, \dots, T$ a large number of samples. As for Gibbs sampling, the order in which we update parameters does not matter. The **R** code for this MH sampler is provided in Panel 1.2.

The reason why in the **R** code we sum the logs of the likelihood and the prior, rather than multiplying the original values, is simply computational. The product of small probabilities can be numbers very close to 0, which computers do not handle well. Thus we add the logarithms, sum, and exponentiate to achieve the desired result. Similarly, in case you have forgotten, $x/y = \exp(\log(x) - \log(y))$, with the latter being favored for computational reasons.

Comparing MH sampling to Gibbs sampling, where all draws from the conditional distribution are used, in the MH algorithm we discard a portion of the candidate values, which inherently makes it less efficient than Gibbs sampling – the price you pay for its increased generality. In Step 1 of the MH sampler we had to choose a variance, δ , for the normal proposal distribution. Choice of the parameters that define our candidate distribution is also referred to as ‘tuning’, and it is important since adequate tuning will make your algorithm more efficient. δ should be chosen (a) large enough so that each step of drawing a new proposal value for θ can cover a reasonable distance in the parameter space, as otherwise, mixing of the Markov chain is inefficient and chains will tend to have strong autocorrelation; and (b) small enough so that proposal values are not rejected too often, as otherwise the random walk will ‘get stuck’ at specific values for too long. As a rule of thumb, your candidate value should be accepted in about 40% of all cases. Acceptance rates of 20 – 80% are probably ok, but anything below or above may well render your algorithm inefficient (this does not mean that it will give you wrong results, only that you will need more iterations to converge to the posterior distribution). In practice, tuning will require some ‘trial-and-error’, some common sense and, with enough experience, some intuition. Or, one can use an adaptive phase, where the tuning parameter is

```

Logreg.MH<-function(y=y, mu0=mu0, sig0=sig0, delta=delta, niter=niter) {

out<-c()

theta<-runif(1, -3,3) #initial value

for (iter in 1:niter){
theta.cand<-rnorm(1, theta, delta)

loglike<-sum(dbinom(y, 1, exp(theta)/(1+exp(theta)), log=TRUE))
logprior <- dnorm(theta,mu0 ,sig0, log=TRUE)
loglike.cand<-sum(dbinom(y, 1, exp(theta.cand)/(1+exp(theta.cand)),
log=TRUE))
logprior.cand <- dnorm(theta.cand, mu0, sig0, log=TRUE)

if (runif(1)<exp((loglike.cand+logprior.cand)-(loglike+logprior))){
theta<-theta.cand
}
out[iter]<-theta
}

return(out)
}

```

Panel 1.2: **R** code to run a Metropolis sampler on a simple logit-normal model.

315 automatically adjusted until it reaches a user-defined acceptance rate, at which
 316 point the adaptive phase ends and the actual Markov chain begins. This is
 317 computationally a little more advanced. Link and Barker (2010) discuss this in
 318 more detail. It is important that the samples drawn during the adaptive phase
 319 are discarded.

320 To illustrate the effects of tuning, we ran the Metropolis-Hastings algorithm
 321 in Panel 1.2 with $\delta = 0.01$, $\delta = 0.2$ and $\delta = 1$. The first 150 iterations for θ are
 322 shown in Fig. 1.3. We see that for a very small δ (the dashed line) the burn-in
 323 is extremely slow - after 150 iterations the chain isn't even half way there, while
 324 for the other two values of δ (solid and dotted) the burn-in phase seems to be
 325 over after only about 10 iterations. While $\delta = 0.2$ leads to reasonably good
 326 mixing, the chain clearly gets stuck on certain values with $\delta = 1$.

327 Other than graphically, you can easily check acceptance rates for the param-
 328 eters you monitor (that are part of your output) using the `rejectionRate()`
 329 function of the package `coda` (we will talk more about this package a little later
 330 on). Do not let the term 'rejection rate' confuse you; it is simply $1 - \text{acceptance}$

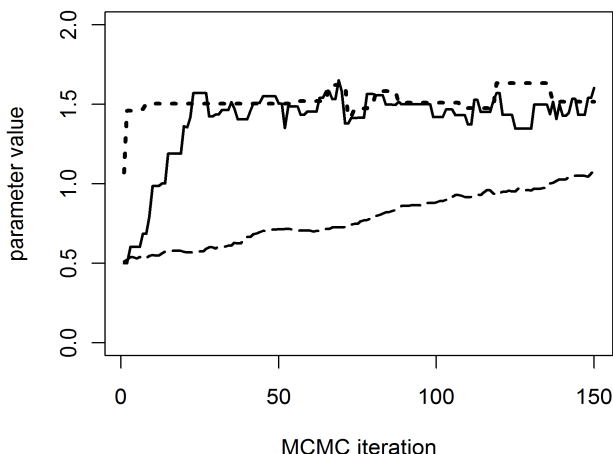


Figure 1.3: Time series plots of θ from a MH algorithm with tuning parameter $\delta = 0.01$ (dashed line), 0.2 (solid line) and 1 (dotted line).

331 rate. There may be parameters – for example, individual values of a random
 332 effect or latent variables – that you do not want to save, though, and in our
 333 next example we will show you a way to monitor their acceptance rates with a
 334 few extra lines of code.

335 1.3.3 Metropolis-within-Gibbs

336 One weakness of the MH sampler is that formulating the joint posterior when
 337 evaluating whether to accept or reject the candidate values for θ becomes in-
 338 creasingly complex or inefficient as the number of parameters in a model in-
 339 creases. As you already saw in Chapt. ??, in these cases you can simply com-
 340 bine MH sampling and Gibbs sampling. You can use the principles of Gibbs
 341 sampling to break down your high-dimensional parameter space into easy-to-
 342 handle one-dimensional conditional distributions and use MH sampling for these
 343 conditional distributions. Better yet, if you have some conjugacy in your model,
 344 you can use the more efficient Gibbs sampling for these parameters and one-
 345 dimensional MH for all the others. You have already seen the basics of how to
 346 build both types of algorithms, so we can jump straight into an example here
 347 and build a Metropolis-within-Gibbs algorithm.

348 **GLMMs: Poisson regression with a random effect** Let's assume a model
 349 that gets us closer to the problem we ultimately want to deal with - a GLMM.
 350 Here, we assume we have Poisson counts, y_{ij} , from $j = 1, 2, \dots, n$ plots in i

different study sites, and we believe that the counts are influenced by some plot-specific covariate, \mathbf{x} , but that there is also a random site effect. So our model is:

$$y_{ij} \sim \text{Poisson}(\lambda_{ij})$$

$$\lambda_{ij} = \exp(\alpha_i + \beta x_{ij})$$

Let's place normal priors on α and β ,

$$\alpha_i \sim \text{Normal}(\mu_\alpha, \sigma_\alpha)$$

and

$$\beta \sim \text{Normal}(\mu_\beta, \sigma_\beta)$$

In this model, we do not specify μ_α and σ_α , but instead, estimate them as well, so we have to specify hyperpriors for these parameters:

$$\begin{aligned} \mu_\alpha &\sim \text{Normal}(\mu_0, \sigma_0) \\ \sigma_\alpha^2 &\sim \text{Inverse-Gamma}(a_0, b_0) \end{aligned}$$

Note that for simplicity we assume that β is constant across the i study sites, and for analysis we set μ_β and σ_β (i.e., we don't estimate these parameters from the data). With the model completely specified, we can compile the full conditionals, breaking the multi-dimensional parameter space into one-dimensional components:

$$\begin{aligned} [\alpha_1 | \alpha_2, \alpha_3, \dots, \alpha_i, \beta, \mathbf{y}_1] &\propto [\mathbf{y}_1 | \alpha_1, \beta][\alpha_1] \\ &\propto \text{Poisson}(\mathbf{y}_1 | \exp(\alpha_1 + \beta \mathbf{x}_1)) \times \text{Normal}(\alpha_1 | \mu_\alpha, \sigma_\alpha), \end{aligned}$$

where $\mathbf{y}_1 = (y_{11}, y_{12}, \dots, y_{1n})$ is the vector of observed counts for site $i = 1$ and, in general, \mathbf{y}_i is the vector of all counts for site i ; analogous, \mathbf{x}_i is the vector of all observations of the covariate for site i . The other full conditionals for each α_i are constructed similarly:

$$\begin{aligned} [\alpha_2 | \alpha_1, \alpha_3, \dots, \alpha_i, \beta, \mathbf{y}_2] &\propto [\mathbf{y}_2 | \alpha_2, \beta][\alpha_2] \\ &\propto \text{Poisson}(\mathbf{y}_2 | \exp(\alpha_2 + \beta \mathbf{x}_2)) \times \text{Normal}(\alpha_2 | \mu_\alpha, \sigma_\alpha), \end{aligned}$$

and so on for all elements of α . The full-conditional for β is:

$$\begin{aligned} [\beta | \alpha, \mathbf{y}] &\propto [\mathbf{y} | \alpha, \beta][\beta] \\ &\propto \text{Poisson}(\mathbf{y} | \exp(\alpha + \beta \mathbf{x})) \times \text{Normal}(\beta | \mu_\beta, \sigma_\beta). \end{aligned}$$

Finally, we need to update the hyperparameters for the random effects vector α :

$$[\mu_\alpha | \alpha] \propto [\alpha | \mu_\alpha, \sigma_\alpha][\mu_\alpha]$$

$$[\sigma_\alpha | \alpha] \propto [\alpha | \mu_\alpha, \sigma_\alpha] [\sigma_\alpha]$$

Note that the likelihood contributions of the counts \mathbf{y} at each site, when conditioned on α , do not depend on the hyperparameters μ_α and σ_α . As such, the full conditionals for these hyperparameters only depend on the collection of all α , not the data. Since we assumed α to come from a normal distribution, the choice of priors for μ_α (normal) and σ_α^2 (inverse-gamma) leads to the same conjugacy we observed in our initial normal model, so that both hyperparameters can be updated using Gibbs sampling.

Now let's build the updating steps for these full conditionals. Again, for the MH steps that update α and β we use normal proposal distributions with standard deviations δ_α and δ_β .

First, we set the initial values $\alpha^{(0)}$ and $\beta^{(0)}$. Then, starting with α_1 , we draw $\alpha_1^{(1)}$ from $\text{Normal}(\alpha_1^{(0)}, \delta_\alpha)$, calculate the conditional posterior density of $\alpha_1^{(0)}$ and $\alpha_1^{(1)}$ and compare their ratios,

$$r = \frac{\text{Poisson}(\mathbf{y}_1 | \exp(\alpha_1^{(1)} + \beta \mathbf{x}_1)) \times \text{Normal}(\alpha_1^{(1)} | \mu_\alpha, \sigma_\alpha)}{\text{Poisson}(\mathbf{y}_1 | \exp(\alpha_1^{(0)} + \beta \mathbf{x}_1)) \times \text{Normal}(\alpha_1^{(0)} | \mu_\alpha, \sigma_\alpha)}$$

and accept $\alpha_1^{(1)}$ with probability $\min(r, 1)$. We repeat this for all α .

For β , we draw $\beta^{(1)}$ from $\text{Normal}(\beta^{(0)}, \delta_\beta)$, compare the posterior densities of $\beta^{(0)}$ and $\beta^{(1)}$,

$$r = \frac{\text{Poisson}(\mathbf{y} | \exp(\alpha + \beta^{(1)} \mathbf{x})) \times \text{Normal}(\beta^{(1)} | \mu_\beta, \sigma_\beta)}{\text{Poisson}(\mathbf{y} | \exp(\alpha + \beta^{(0)} \mathbf{x})) \times \text{Normal}(\beta^{(0)} | \mu_\beta, \sigma_\beta)},$$

and accept $\beta^{(1)}$ with probability $\min(r, 1)$.

For μ_α and σ_α^2 , we sample directly from the full conditional distributions (Eq. 1.2 and Eq. 1.3):

$$\mu_\alpha^{(1)} \sim \text{Normal}(\mu_n, \sigma_n^2)$$

where

$$\mu_n = \frac{\sigma_\alpha^{2(0)}}{\sigma_\alpha^{2(0)} + n_\alpha \sigma_0^2} \times \mu_0 + \frac{n_\alpha \sigma_0^2}{\sigma_\alpha^{2(0)} + n_\alpha \sigma_0^2} \times \bar{\alpha}^{(1)}$$

and

$$\sigma_n^2 = \frac{\sigma_\alpha^{2(0)} \sigma_0}{\sigma_\alpha^{2(0)} + n \sigma_0^2}$$

Here, $\bar{\alpha}$ is the current mean of the vector α , which we updated before, and n_α is the length of α . For σ_α^2 we use

$$\sigma_\alpha^{2(1)} \sim \text{Inverse-Gamma}(a_n, b_n),$$

where

$$a_n = n_a/2 + a_0,$$

396 and

$$b_n = 0.5 \sum_{i=1}^{n_\alpha} (\alpha_i^{(1)} - \mu_\alpha^{(1)})^2 + b_0.$$

397 We repeat these steps over T iterations of the MCMC algorithm. Call the
 398 function `PoisGLMM()` in `scrbook` to check out what this algorithm looks like in
 399 **R**.

400 In this example we may not want to save each individual α_i , but are only
 401 interested in their mean and standard deviation. Since these two parameters will
 402 change as soon as the value for one element in α changes, their acceptance rates
 403 will always be close to 1 and are not representative of how well your algorithm
 404 performs. To monitor the acceptance rates of parameters you do not want to
 405 save, you simply need to add a few lines of code into your updater to see how
 406 often the individual parameters are accepted. The code for updating α from our
 407 Poisson GLMM below shows one way how to monitor acceptance of individual
 408 α_i 's.

```

409 #initiate counter for acceptance rate of alpha
410 alphaUps<-0
411
412 #loop over sites, update intercepts alpha one at a time;
413 #only data at site i contributes information
414 #lev is the number of sites i
415 for (i in 1:lev) {
416   alpha.cand<-rnorm(1, alpha[i], delta_alpha)
417   loglike<- sum(dpois (y[site==i], exp(alpha[i] + beta*x[site==i]),
418     log=TRUE))
419   logprior<- dnorm(alpha[i], mu_alpha,sig_alpha, log=TRUE)
420   loglike.cand<- sum(dpois (y[site==i], exp(alpha.cand + beta *x[site==i]),
421     log=TRUE))
422   logprior.cand<- dnorm(alpha.cand, mu_alpha,sig_alpha, log=TRUE)
423   if (runif(1)< exp((loglike.cand+logprior.cand) -(loglike+logprior))) {
424     alpha[i]<-alpha.cand
425     alphaUps<-alphaUps+1
426   }
427 }
428
429 #lets you check the acceptance rate of alpha at every 100th iteration
430 if(iter %% 100 == 0) {
431   cat("    Acceptance rates\n")
432   cat("    alpha =", alphaUps/lev, "\n")
433 }
```

434 1.3.4 Rejection sampling and slice sampling

435 While MH and Gibbs sampling are probably the most widely applied algorithms
 436 for posterior approximation, there are other options that work under certain
 437 circumstances and may be more efficient when applicable. **WinBUGS** applies

these algorithms and we want you to be aware that there is more out there to approximate posterior distributions than Gibbs and MH. One alternative algorithm is rejection sampling. Rejection sampling is not an MCMC method, since each draw is independent of the others. The method can be used when the posterior $[\theta|y]$ is not a known parametric distribution but can be expressed in closed form. Then, we can use a so-called envelope function, say, $g(\theta)$, that we can easily sample from, with the restriction that $[\theta|y] < M \times g(\theta)$. We then sample a candidate value for θ from $g(\theta)$, calculate $r = [\theta|y]/M \times g(\theta)$ and keep the sample with the probability r . M is a constant that has to be picked so that r lies between 0 and 1, for example by evaluating both $[\theta|y]$ and $g(\theta)$ at n points and looking at their ratios. Rejection sampling only works well if $g(\theta)$ is similar to $[\theta|y]$, and packages like **WinBUGS** use adaptive rejection sampling (Gilks and Wild, 1992), where a complex algorithm is used to fit an adequate and efficient $g(\theta)$ based on the first few draws. Though efficient in some situations, rejection sampling does not work well with high-dimensional problems, since it becomes increasingly hard to define a reasonable envelope function. For an example of rejection sampling in the context of SCR models, see Chapt. ??, where we use it to simulation inhomogeneous point processes.

Another alternative is slice sampling (Neal, 2003). In slice sampling, we sample uniformly from the area under the plot of $[\theta|y]$. Considering a single univariate θ . Let's define an auxiliary variable, $U \sim \text{Unif}(0, [\theta|y])$. Then, θ can be sampled from the vertical slice of $[\theta|y]$ at U (Fig. 1.4):

$$\theta|U \sim \text{Unif}(B),$$

where $B = \{\theta : [\theta|y] \geq U\}$

Slice sampling can be applied in many situations; however, implementing an efficient slice sampling procedure can be complicated. We refer the interested reader to Robert and Casella (2010, Chapt. 7) for a simple example. Both rejection sampling and slice sampling can be applied on one-dimensional conditional distributions within a Gibbs sampling setup.

1.4 MCMC for Closed Capture-Recapture Model M_h

By now you have seen MCMC samplers for some simple generalized (mixed) linear models. Now, to ease you into more complex models, we construct our own MCMC algorithm using a Metropolis-within-Gibbs sampler for the non-spatial model with individual heterogeneity in capture probability, model M_h , developed in Chapt. ??.

To recapitulate: Under the non-spatial model, each of the n observed individuals is either detected (1) or not (0) during each of K sampling occasions. We estimate N using data augmentation and have a Bernoulli model for the data augmentation variables z_i .

$$z_i \sim \text{Bernoulli}(\psi)$$

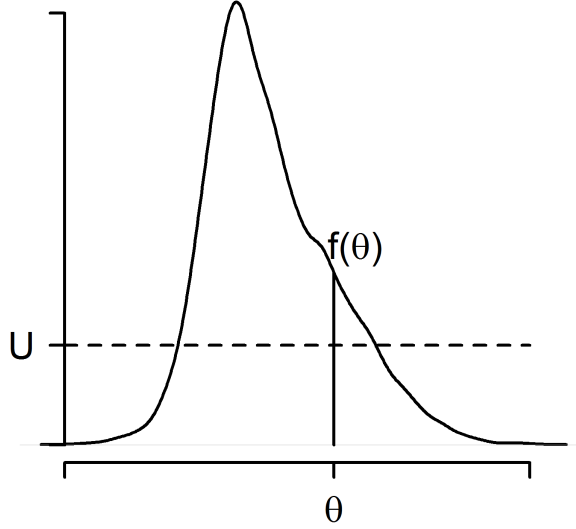


Figure 1.4: Slice sampling. For $U \sim \text{Unif}(0, [\theta|y])$, we can sample θ from the vertical slice of $[\theta|y]$ at U ; $\theta|U \sim \text{Unif}(B)$, where $B = \{\theta : [\theta|y] \geq U\}$.

477 The binomial observation model is expressed conditional on the latent variables
 478 z_i .

$$y_i \sim \text{Binomial}(p_i \times z_i, K)$$

479 Further, we prescribe a distribution for the capture probability p_i . Here we
 480 assume

$$\text{logit}(p_i) \sim \text{Normal}(\mu_p, \sigma_p^2)$$

481 As usual, we have to go through two general steps before we write the MCMC
 482 algorithm:

- 483 (1) Identify the model with all its components (including priors)
- 484 (2) Recognize and express the full conditional distributions for all parameters

485 Our model components are as follows: $[y_i|p_i, z_i]$, $[p_i|\mu_p, \sigma_p]$, and $[z_i|\psi]$ for *each*
 486 $i = 1, 2, \dots, M$ and then prior distributions $[\mu_p]$, $[\sigma_p]$ and $[\psi]$. The joint posterior
 487 distribution of all unknown quantities in the model is proportional to the joint
 488 distribution of all elements y_i, p_i, z_i and also the prior distributions of the prior
 489 parameters:

$$\left\{ \prod_{i=1}^M [y_i|p_i, z_i][p_i|\mu_p, \sigma_p][z_i|\psi] \right\} [\mu_p, \sigma_p, \psi]$$

490 For prior distributions, we assume that μ_p, σ_p, ψ are mutually independent and
 491 for μ_p and σ_p we use improper uniform priors, and $\psi \sim \text{Uniform}(0, 1)$. This is

equivalent to $\text{Beta}(1, 1)$, which will come in handy, as we will see in a moment. Note that the likelihood contribution for each individual, when conditioned on p_i and z_i , does not depend on ψ , μ_p , or σ_p . As such, the full-conditional for the structural parameter ψ only depend on the collection of data augmentation variables z_i , and that for μ_p and σ_p will only depend on the collection of latent variables $p_i; i = 1, 2, \dots, M$ (this is equivalent to what we saw in the Poisson regression with random intercept α , where hyperparameters for the distribution of α did not depend on the observed data). The full conditionals for all the unknowns are as follows:

(1) For p_i :

$$[p_i | y_i, \mu_p, \sigma_p, z_i] \propto \begin{cases} [y_i | p_i][p_i | \mu_p, \sigma_p] & \text{if } z_i = 1 \\ [p_i | \mu_p, \sigma_p] & \text{if } z_i = 0 \end{cases}$$

(2) for z_i :

$$[z_i | y_i, p_i, \psi] \propto [y_i | z_i \times p_i] \text{Bernoulli}(z_i | \psi)$$

(3) For μ_p :

$$[\mu_p | p_i, \sigma_p] \sim \left\{ \prod_i [p_i | \mu_p, \sigma_p] \right\} \times \text{const}$$

(4) For σ_p :

$$[\sigma_p | p_i, \mu_p] \sim \left\{ \prod_i [p_i | \mu_p, \sigma_p] \right\} \times \text{const}$$

(5) For ψ :

$$[\psi | z_i] \propto \left\{ \prod_i [z_i | \psi] \right\} \text{Beta}(1, 1)$$

Remember that $\text{Beta}(1, 1)$ is equivalent to $\text{Uniform}(0, 1)$. The beta distribution is the conjugate prior to the binomial and Bernoulli distributions and the general form of a full conditional of a beta-binomial model with $x_i \sim \text{Bernoulli}(\theta)$ and $\theta \sim \text{Beta}(a, b)$ is

$$[\theta | \mathbf{x}] \propto \text{Beta}(a + \sum_i x_i, b + n - \sum_i x_i).$$

In our case that means

$$[\psi | z_i] \propto \text{Beta}(1 + \sum z_i, 1 + M - \sum z_i).$$

What we've done here is identify each of the full conditional distributions in sufficient detail to toss them into our Metropolis-Hastings algorithm. The constant terms in the full conditionals for μ_p and σ_p reflect the improper prior we chose for both parameters. Because of the choice of an improper prior, prior probability densities for both parameters $\propto 1$, i.e. constant, and these constants cancel out of the MH acceptance ratio (see updating step below and following

example). Below, you see the updating step for the detection parameter \mathbf{p} . Note that (1) we draw candidate values on the logit scale and (2) instead of looping through $1 - M$ individuals to update all p_i , we update all elements of the vector of \mathbf{p} in parallel, for computational efficiency.

```

521 ### update the logit(p) parameters
522 lp.cand<- rnorm(M,lp,1) # 1 is a tuning parameter
523 p.cand<-plogis(lp.cand)
524 ll<-dbinom(ytot,K,z*p, log=T)
525 prior<-dnorm(lp,mu,sigma, log=T)
526 llcand<-dbinom(ytot,K,z*p.cand, log=T)
527 prior.cand<-dnorm(lp.cand,mu,sigma, log=T)
528
529 kp<- runif(M) < exp((llcand+prior.cand)-(ll+prior))
530 p[kp]<-p.cand[kp]
531 lp[kp]<-lp.cand[kp]

```

The parameters μ_p and σ_p are also updated using MH steps (see the code for μ_p below). In truth, we could also sample μ_p and σ_p^2 directly with certain choices of prior distributions. For example, if $\mu_p \sim \text{Normal}(0,1000)$ then the full conditional for μ_p is also normal (see sec. 1.3.1), etc..

```

536 p0.cand<- rnorm(1,p0,.05)
537 if(p0.cand>0 & p0.cand<1){
538   mu.cand<-log(p0.cand/(1-p0.cand))
539   ll<-sum(dnorm(lp,mu,sigma,log=TRUE))
540   llcand<-sum(dnorm(lp,mu.cand,sigma,log=TRUE))
541   if(runif(1)<exp(llcand-ll)) {
542     mu<-mu.cand
543     p0<-p0.cand
544   }
545 }

```

For ψ we can easily sample directly from the beta distribution:

```

547 psi<-rbeta(1, sum(z) + 1, M-sum(z) + 1)

```

To update the z_i we have opted for a MH updater (although they could be updated directly from their full-conditional). Since z_i can only take the values of 0 or 1, we generate candidate values using $\mathbf{z.cand}<-ifelse(\mathbf{z}=1,0,1)$. The updating step for z_i is detailed in the next example. You can check out the full code by invoking `modelMh()` from the **R** package `scrbook`.

553 1.5 MCMC Algorithm for Model SCR0

Conceptually, but also in terms of MCMC coding, it is only a small step from the non-spatial model M_h to a fully spatial capture-recapture model. Next, we

will walk you through the steps of building your own MCMC sampler for the basic SCR model (i.e. without any individual, site or time specific covariates) with both a Poisson and a binomial encounter process. As usual, we will have to go through two general steps before we write the MCMC algorithm:

- (1) Identify the model with all its components (including priors)
- (2) Recognize and express the full conditional distributions for all parameters

It is worthwhile to go through all of step 1 for an SCR model, but you have probably seen enough of step 2 in our previous examples to get the essence of how to express a full conditional distribution. Therefore, we will exemplify step 2 for some parameters and tie these examples directly to the respective R code.

Step 1 – Identify your model Recall the components of the basic SCR model with a Poisson encounter process from Chapt. ??: We assume that individuals i , or rather, their activity centers \mathbf{s}_i , are uniformly distributed across the state-space \mathcal{S} ,

$$\mathbf{s}_i \sim \text{Uniform}(\mathcal{S})$$

and that the number of times individual i encounters trap j , y_{ij} , is a Poisson variable with mean λ_{ij} ,

$$y_{ij} \sim \text{Poisson}(\lambda_{ij}).$$

The link between individual location, movement and trap encounter rates is made by the assumption that λ_{ij} , is a decreasing function of the distance between \mathbf{s}_i and the location of j , \mathbf{x}_j , say

$$d_{ij} = \|\mathbf{s}_i - \mathbf{x}_j\|,$$

of the Gaussian (or half-normal) form

$$\lambda_{ij} = \lambda_0 \exp(-d_{ij}^2/2\sigma^2),$$

where λ_0 is the baseline trap encounter rate at $d_{ij} = 0$ and σ is the scale parameter of the half-normal function.

As in the non-spatial exaple for model M_h , we estimate N , here the number of \mathbf{s}_i in \mathcal{S} , using data augmentation (sec. ??). We create $M-n$ all-zero encounter histories and estimate N by summing over the auxiliary data augmentation variables, z_i , which we assume is a Bernoulli random variable,

$$z_i \sim \text{Bernoulli}(\psi).$$

To link the two model components, we modify our trap encounter model to

$$\lambda_{ij} = \lambda_0 \times \exp(-d_{ij}^2/2\sigma^2) \times z_i.$$

The model has the following structural parameters, for which we need to specify priors:

585 ψ : the Uniform(0, 1) is required as part of the data augmentation procedure
 586 and in general is a natural choice of an uninformative prior for a proba-
 587 bility. It will also lead to conjugacy as we saw in the example of model
 588 M_h , so that we can update ψ directly from its full conditional distribution
 589 using Gibbs sampling.

590 \mathbf{s}_i : since \mathbf{s}_i is a pair of coordinates it is two-dimensional and we use a uniform
 591 prior limited by the extent of our state-space over both dimensions.

592 σ : we can conceive several priors for σ but let's assume an improper prior,
 593 one that is Uniform over $(0, \infty)$. As we already saw, this choice is conve-
 594 nient when updating the parameter, because the constant prior probability
 595 cancels out of the MH acceptance ratio.

596 λ_0 : analogous, we will use a Uniform(0, ∞) improper prior for λ_0 .

597 **Step 2 – Construct the full conditionals:** Having completed step 1,
 598 let's look at the full conditional distributions for some of these parameters. We
 599 saw that with improper priors, full conditionals are proportional only to the
 600 likelihood of the observations; for example, consider σ :

$$[\sigma | \mathbf{s}, \lambda_0, \mathbf{z}, \mathbf{y}] \propto \left\{ \prod_i [y_i | \mathbf{s}_i, \lambda_0, z_i, \sigma] \right\}$$

601 The **R** code to update σ is shown below. Notice that we automatically reject
 602 negative candidate values, since σ cannot be < 0 .

```
603 sig.cand <- rnorm(1, sigma, 0.1) #draw candidate value
604 if(sig.cand>0){ #automatically reject sig.cand that are <0
605   lam.cand <- lam0*exp(-(d*d)/(2*sig.cand*sig.cand))
606   ll<- sum(dpois(y, lam*z, log=TRUE))
607   llcand <- sum(dpois(y, lam.cand*z, log=TRUE))
608   if(runif(1) < exp( llcand - ll) ){
609     ll<-llcand
610     lam<-lam.cand
611     sigma<-sig.cand
612   }
613 }
```

614 These steps are analogous for λ_0 and \mathbf{s}_i and we will use MH steps for all of
 615 these parameters. Similar to the random intercepts in our Poisson GLMM, we
 616 update each \mathbf{s}_i individually. Note that to be fully correct, the full conditional
 617 for \mathbf{s}_i contains both the likelihood and prior component, since we did not specify
 618 an improper, but a proper uniform prior on \mathbf{s}_i . However, with a uniform dis-
 619 tribution the probability density of any value is $1/(\text{upper limit} - \text{lower limit}) =$
 620 constant. Thus, the prior components are identical for both the current and the
 621 candidate value so that when you calculate the ratio of posterior densities, r ,

the identical prior component appears both in the numerator and denominator and cancel each other out.

We still have to update z_i . The full conditional for z_i is

$$[z_i | y_i, \sigma, \lambda_0, \mathbf{s}_i] \propto [y_i | z_i, \sigma, \lambda_0, \mathbf{s}_i][z_i]$$

and since $z_i \sim \text{Bernoulli}(\psi)$, the term has to be taken into account when updating z_i :

```

627     zUps <- 0 #set counter to monitor acceptance rate
628     for(i in 1:M) {
629       #no need to update seen individuals, since their z =1
630       if(seen[i])
631         next
632       zcand <- ifelse(z[i]==0, 1, 0)
633       llz <- sum(dpois(y[i,], lam[i,]*z[i], log=TRUE))
634       llcand <- sum(dpois(y[i,], lam[i,]*zcand, log=TRUE))
635
636       prior <- dbinom(z[i], 1, psi, log=TRUE)
637       prior.cand <- dbinom(zcand, 1, psi, log=TRUE)
638       if(runif(1) < exp((llcand+prior.cand)-(llz+prior))){
639         z[i] <- zcand
640         zUps <- zUps+1
641       }
642     }

```

The parameter ψ is a hyperparameter of the model, with an uninformative prior distribution of Uniform(0,1) or Beta(1,1), so that

$$[\psi | \mathbf{z}] \propto \text{Beta}(1 + \sum_i z_i, 1 + M - \sum_i z_i).$$

These are all the building blocks you need to write the MCMC algorithm for the spatial null model with a Poisson encounter process. You can find the full **R** code by calling the function (**SCR0pois**) in the **R** package **scrbook**.

1.5.1 SCR model with binomial encounter process

The equivalent SCR model with a binomial encounter process is very similar. Here, each individual i can only be detected once at any given trap j during a sampling occasion k . Thus

$$y_{ij} \sim \text{Binomial}(p_{ij}, K)$$

Where p_{ij} is some function of distance between \mathbf{s}_i and trap location \mathbf{x}_j . Here we use:

$$p_{ij} = 1 - \exp(-\lambda_{ij})$$

Recall from Chapt. ?? that this is the complementary log-log (cloglog) link function, which constrains p_{ij} to fall between 0 and 1. For our MCMC algorithm that means that, instead of using a Poisson likelihood, $\text{Poisson}(y|\sigma, \lambda_0, \mathbf{s}, z)$, we use a binomial likelihood, $\text{Binomial}(y|\sigma, \lambda_0, \mathbf{s}, z; K)$, in all the conditional distributions. An exemplary updating step for λ_0 under a binomial encounter model is shown below. The full MCMC code for the binomial SCR with a cloglog link (`SCR0binom.cl`) can be found in the **R** package `scrbook`.

```

661     lam0.cand <- rnorm(1, lam0, 0.1)
662     #automatically reject lam0.cand that are <0
663     if(lam0.cand >0){
664         lam.cand <- lam0.cand*exp(-(d*d)/(2*sigma*sigma))
665         p.cand <- 1-exp(-lam.cand)
666         ll<- sum(dbinom(y, K, pmat *z, log=TRUE))
667         llcand <- sum(dbinom(y, K, p.cand *z, log=TRUE))
668         if(runif(1) < exp( llcand - ll )){
669             ll<-llcand
670             pmat<-p.cand
671             lam0<- lam0.cand
672         }
673     }

```

Another possibility is to model variation in the individual and site specific detection probability, p_{ij} , directly, without any transformation, such that

$$p_{ij} = p_0 \times \exp(-d_{ij}^2/(2\sigma^2))$$

and $p_0 \in [0, 1]$. This formulation is analogous to how detection probability is modeled in distance sampling under a half-normal detection function; however, in distance sampling p_0 – detection of an individual on the transect line – is assumed to be 1 (Buckland et al., 2001). Under this formulation the updater for p_0 becomes:

```

681     p0.cand <- rnorm(1, p0, 0.1)
682     if(p0.cand >0 & p0.cand < 1 ){
683         #automatically rejects lam0.cand that are not {0,1}
684         p.cand <- p0.cand*exp(-(d*d)/(2*sigma*sigma))
685         ll<- sum(dbinom(y, K, pmat *z, log=TRUE))
686         llcand <- sum(dbinom(y, K, p.cand *z, log=TRUE))
687         if(runif(1) < exp( llcand - ll )){
688             ll<-llcand
689             pmat<-p.cand
690             p0<- p0.cand
691         }
692     }

```


693 1.6 Looking at Model Output

694 Now that you have an MCMC algorithm to analyze spatial capture-recapture
 695 data with, let's run an actual analysis so we can look at the output. As an
 696 example, we will use the Fort Drum bear data set we first introduced in Chapt.
 697 ?? and already analyzed in several preceding chapters. You can load the Fort
 698 Drum data (`data(beardata)`), extract the trap locations (`trapmat`) and de-
 699 tection data (`bearArray`) and build the augmented $M \times J$ array of individual
 700 encounter histories:

```
701 > M=700
702 > trapmat<-beardata$trapmat
703 #summarizes captures across occasions
704 > bearmat<-apply(beardata$bearArray, 1:2, sum)
705 > Xaug<-matrix(0, nrow=M, ncol=dim(trapmat)[1])
706 > Xaug[1:dim(bearmat)[1],]<-bearmat #create augmented data set
```

707 In addition to these data, we need to specify the outermost coordinates of
 708 the state-space. Since bears are wide ranging animals we add a 20-km buffer
 709 to the maximum and minimum coordinates of the trap array:

```
710 > xl<- min(trapmat[,1])- 20
711 > yl<- min(trapmat[,2])- 20
712 > xu<- max(trapmat[,1])+ 20
713 > yu<- max(trapmat[,2])+ 20
```

714 Finally, use the MCMC code for the binomial encounter model with the
 715 cloglog link (`SCR0binom.cl`) and run 5000 iterations. This should take approxi-
 716 mately 25 minutes (in real life we would of course run the algorithm a lot longer
 717 but for demonstration purposes let's stick with a number of iterations that can
 718 be run in a manageable amount of time).

```
719 > set.seed(13)
720 > mod0<-SCR0binom.cl(y=Xaug, X=trapmat, M=M, xl=xl, xu=xu, yl=yl,
721 +                    yu=yu, K=8, delta=c(0.1, 0.05, 2), niter=5000)
```

722 Before, we used simple **R** commands to look at model results. However, there
 723 is a specific **R** package to summarize MCMC simulation output and perform
 724 some convergence diagnostics – package `coda` (Plummer et al., 2006). Download
 725 and install `coda`, then convert your model output to an `mcmc` object

```
726 > chain<-mcmc(mod0)
```

727 which can be used by `coda` to produce MCMC specific output.

728 1.6.1 Markov chain time series plots

729 Start by looking at time series plots of your Markov chains using `plot(chain)`.
 730 This command produces a time series plot and marginal posterior density plots

for each monitored parameter, similar to what we did before using the `hist()` and `plot()` commands. Fig. 1.5 shows an example of these plots for σ and λ_0 . Time series plots will tell you several things: First, recall from sec. 1.3.2 that the way the chains move through the parameter space gives you an idea of whether your MH steps are well tuned. If chains were constant over many iterations you would need to decrease the tuning parameter of the (normal) proposal distribution. If a chain moves along some gradient to a stationary state very slowly, you may want to increase the tuning parameter so that the parameter space is explored more efficiently.

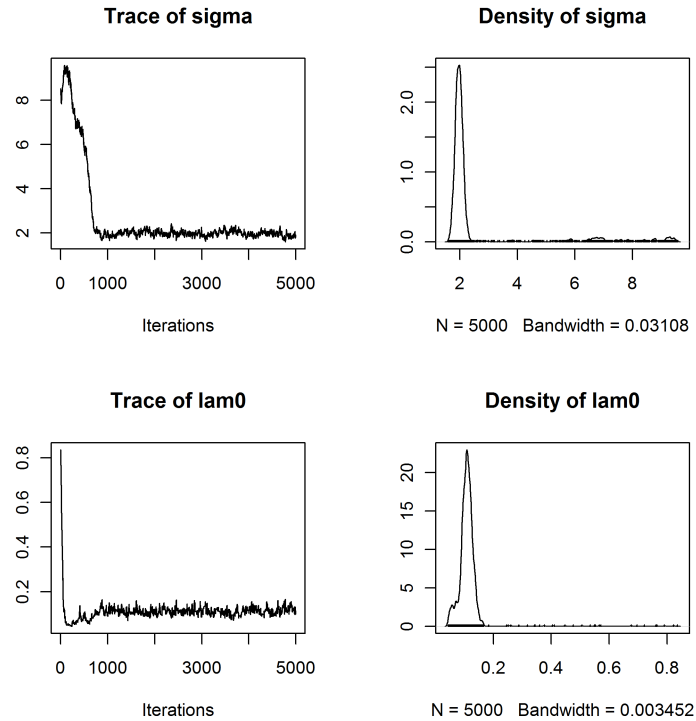


Figure 1.5: Time series and posterior density plots for σ and λ_0 for the Fort Drum black bear data.

Second, you will be able to see if your chains converged and how many initial simulations you have to discard as burn-in. In the case of the chains shown in Fig. 1.5, we would probably consider the first 750 – 1000 iterations as burn-in, as afterwards the chains seem to be fairly stationary.

1.6.2 Posterior density plots

The `plot()` command also produces posterior density plots and it is worthwhile to look at those carefully. For parameters with priors that have bounds (e.g. uniform over some interval), you will be able to see if your choice of the prior is truncating the posterior distribution. In the context of SCR models, this will mostly involve our choice of M , the size of the augmented data set. If the posterior of N has a lot of mass concentrated close to M (or equivalently the posterior of ψ has a lot of mass concentrated close to 1), as in the example in Fig. 1.6, we have to re-run the analysis with a larger M . A diffuse posterior plot suggests that the parameter may not be well-identified. There may not be enough information in your data to estimate model parameters and you may have to consider a simpler model. Finally, posterior density plots will show you if the posterior distribution is symmetrical or skewed – if the distribution has a heavy tail, using the mean as a point estimate of your parameter of interest may be biased and you may want to opt for the median or mode instead.

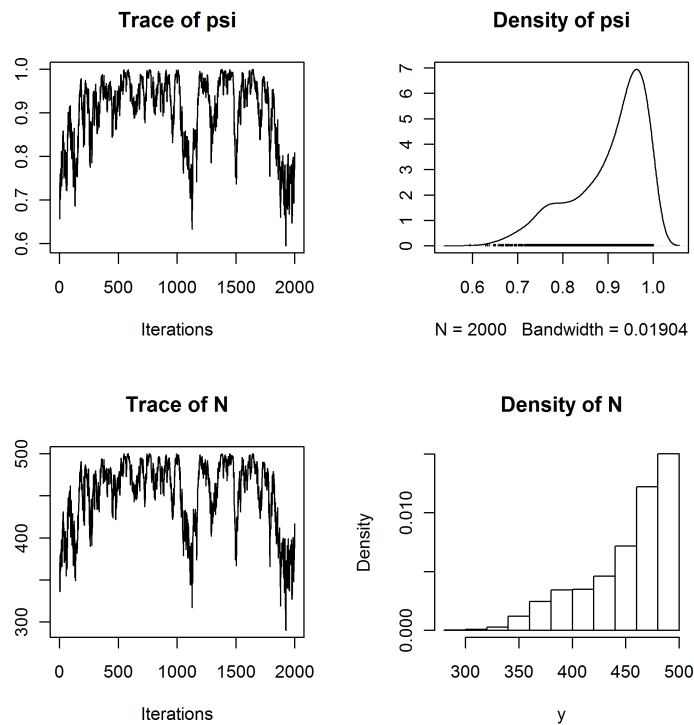


Figure 1.6: Time series and posterior density plots of ψ and N for the Fort Drum black bear data truncated by the upper limit of M (500).

1.6.3 Serial autocorrelation and effective sample size

Checking the degree of autocorrelation in your Markov chains and estimating the effective sample size your chain has generated should be part of evaluating your model output. If you use **WinBUGS** through the **R2WinBUGS** package, the `print()` command will automatically return the effective sample size for all monitored parameters. In the **coda** package there are several functions you can use to do so. The function `effectiveSize()` will directly give you an estimate of the effective sample size for the parameters:

```
> effectiveSize(window(chain, start=1001))
      sigma      lam0      psi      N
93.89807 163.72311  51.96443  46.45394
```

Alternatively, you can use the `autocorr.diag()` function, which will show you the degree of autocorrelation for different lag values (which you can specify within the function call, we use the defaults below):

```
> autocorr.diag(window(chain, start=1001))
      sigma      lam0      psi      N
Lag 0  1.0000000 1.0000000 1.0000000 1.0000000
Lag 1  0.9316928 0.91464875 0.9745833 0.9663320
Lag 5  0.7603332 0.67445407 0.8525272 0.8500215
Lag 10 0.6065374 0.48724122 0.7514657 0.7530124
Lag 50 0.1122331 0.06564406 0.3811939 0.3823236
```

In the present case we see that autocorrelation is especially high for the parameter ψ and effective sample size for this parameter is only 52! This means we would have to run the model for much longer to obtain a reasonable effective sample size. Unfortunately, with many SCR data sets we observe high degrees of serial autocorrelation. For now, let's continue using this small number of samples to look at the output.

1.6.4 Summary results

Now that we checked that our chains apparently have converged and pretending that we have generated enough samples from the posterior distribution, we can look at the actual parameter estimates. The `summary()` function will return two sets of results: the mean parameter estimates, with their standard deviation, the naïve standard error – i.e. your regular standard error calculated for T (= number of iterations) samples without accounting for serial autocorrelation – and the Time-series SE (in **WinBUGS** and earlier in this book referred to as MC error), which accounts for autocorrelation. Remember our rule of thumb that this error decreases with increasing chain length and should be 1% or less of the parameter estimate. In **WinBUGS** the MC error is only given in the log output within **BUGS** itself. You should adjust the `summary()` call by removing the burn-in from calculating parameter summary statistics. To do so, use the `window()` command, which lets you specify at which iteration to start

'counting'. In contrast to **WinBUGS**, which requires you to set the burn-in length before you run the model, this command gives us full flexibility to make decisions about the burn-in after we have seen the trajectories of our Markov chains. For our example, `summary(window(chain, start=1001))` returns the following output:

```

805 Iterations = 1001:5000
806 Thinning interval = 1
807 Number of chains = 1
808 Sample size per chain = 4000
809
810 1. Empirical mean and standard deviation for each variable,
811    plus standard error of the mean:
812
813           Mean          SD Naive SE Time-series SE
814 sigma    1.9697    0.12534 0.0019818      0.012792
815 lam0     0.1124    0.01521 0.0002405      0.001311
816 psi      0.7295    0.11794 0.0018648      0.015278
817 N       510.9190  81.99868 1.2965130      10.580567
818
819 2. Quantiles for each variable:
820
821           2.5%        25%         50%         75%        97.5%
822 sigma    1.7288     1.8831     1.9666     2.0517     2.2240
823 lam0     0.0863     0.1008     0.1112     0.1217     0.1449
824 psi      0.5100     0.6423     0.7261     0.8170     0.9549
825 N       359.0000  451.0000  508.0000  572.0000  668.0000

```

Looking at the MC errors (column labeled **Time-series SE**), we see that in spite of the high autocorrelation, the MC error for σ is below the 1% threshold, whereas for all other parameters, MC errors are still above, another indication that for a thorough analysis we should run a longer chain.

Our algorithm gives us a posterior distribution of N , but we are usually interested in the density, D . Density itself is not a parameter of our model, but we can derive a posterior distribution for D by dividing each value of N (N at each iteration) by the area of the state-space (here 3032.719 km²) and we can use summary statistics of the resulting distribution to characterize D :

```

835 > summary(window(chain[,4]/ 3032.719, start=1001))
836
837 Iterations = 1001:5000
838 Thinning interval = 1
839 Number of chains = 1
840 Sample size per chain = 4000
841
842 1. Empirical mean and standard deviation for each variable,
843    plus standard error of the mean:
844

```

```

845           Mean           SD           Naive SE Time-series SE
846           0.1684690       0.0270380       0.0004275       0.0034888

```

```

847
848 2. Quantiles for each variable:

```

```

849
850      2.5%    25%    50%    75%  97.5%
851 0.1184 0.1487 0.1675 0.1886 0.2203

```

```

852 We see that our mean density of  $0.17/km^2$  is very similar to the estimate of
853  $0.18/km^2$  obtained under the non-spatial model  $M_0$  in Chapt. ??.
```

854 1.6.5 Other useful commands

```

855 While inspecting the time series plot gives you a first idea of how well you
856 tuned your MH algorithm, use rejectionRate() to obtain the rejection rates
857 ( $1 - \text{acceptance rates}$ ) of the parameters that are written to your output:

```

```

858 > rejectionRate(chain)
859      sigma      lam0      psi      N
860 0.42988598 0.78775755 0.00000000 0.03160632

```

```

861 Recall (sec. 1.3.2) that rejection rates should lie between 0.2 and 0.8, so our
862 tuning seems to have been appropriate here. Draws of the parameter  $\psi$  are never
863 rejected since we update it with Gibbs sampling, where all candidate values are
864 kept. And since  $N$  is the sum of all  $z_i$ , all it takes for  $N$  to change from one
865 iteration to the next are small changes in the  $z$ -vector, so the rejection rate of  $N$ 
866 is always low. If you have run several parallel chains, you can combine them into
867 a single mcmc object using the mcmc.list() command on the individual chains
868 (note that each chain has to be converted to an mcmc object before combining
869 them with mcmc.list()). You can then easily obtain the Gelman-Rubin diag-
870 nostic (Gelman et al., 2004), in WinBUGS called Rhat, using gelman.diag(),
871 which will indicate if all chains have converged to the same stationary distribu-
872 tion. For details on these and other functions, see the coda manual, which can
873 be found (together with the package) on the CRAN mirror.
```

874 1.7 Manipulating the State-Space

```

875 So far, we have constrained the location of the activity centers to fall within
876 the outermost coordinates of our rectangular state-space by posing upper and
877 lower bounds for  $x$  and  $y$ . But what if  $\mathcal{S}$  has an irregular shape – maybe there
878 is a large water body we would like to remove from  $\mathcal{S}$ , because we know our
879 terrestrial study species does not occur there. Or the study takes place in a
880 clearly defined area such as an island.
```

```

881     As mentioned before, this situation is difficult to handle in BUGS engines.
882 In some simple cases we can adjust the state-space by setting one of the coor-
883 dinates of  $\mathbf{s}_i$  to be some function of the other and reject candidate  $\mathbf{s}_i$  that do

```

not fall within this modified state-space. In this manner, we can cut off corners of the rectangle to approximate the actual state-space³. To visualize this approach, plot the following rectangle, representing your state-space polygon, and line, representing, for example, the approximation of a shore line:

```

> xlim<-c(-5,5)
> ylim<-c(-7,7)
> plot(xlim, ylim, type='n')
> abline(a=4, b=0.4)

```

The Y coordinates limiting your state-space to the habitat that is suitable to the species you study can now be expressed as a linear function of the X coordinates, in this case, $Y = 4 + 0.4 \times X$. To include this new limit in a **BUGS** model, we need to change the following:

```

#draw SX and SY as before
SX[i]~dunif(xlim[1],xlim[2])
SY[i]~dunif(ylim[1],ylim[2])
#calculate upper limit for Y given X
ymax[i]<-4+0.4*SX[i]
# use step function to see if location [SX, SY]
# is below the Y limit (Pin = 1) or not (Pin = 0)
Pin[i] <- step(ymax[i] - SY[i])
In[i] ~ dbern(Pin[i])

```

The object **In** is a vector of M 1's, passed as data to the model. If $\text{Pin} = 0$, the likelihood will be 0 and the candidate $[SX, SY]$ pair will be rejected. If $\text{Pin} = 1$, this bit of the likelihood is equal to 1, and whether or not the the candidate pair of coordinates is accepted depends only on capture history of i . This approach can be very useful in some situations but is clearly restricted by the functional form of the relationship between SX and SY that it requires.

In **R**, we are much more flexible, as we can use the actual state-space polygon to constrain s_i . To illustrate that, let's look at a camera trapping study of raccoons (*Procyon lotor*) conducted on South Core Banks, a barrier island within Cape Lookout National Seashore, North Carolina (details of the study can be found in Sollmann et al. (2013) and in Chapt. ?? where we present the analysis of this data set with spatial mark-resight models). Since camera-traps were spread across the entire length of the island, we set the state-space to be delineated by the shore line of the island (Fig. 1.7), which clearly cannot easily be approximated as a rectangle. Instead, within **R** we can use an actual shapefile of the island.

In other circumstances you may still want to create the state-space as before, by adding some buffer to your trapping grid, but you may find that the resulting rectangle includes water bodies, paved parking lots or any other kind of habitat you know is never used by the species you study. In order to precisely describe

³This idea was pitched to us by Mike Meredith, Biodiversity Conservation Society Sarawak/WCS Malaysia

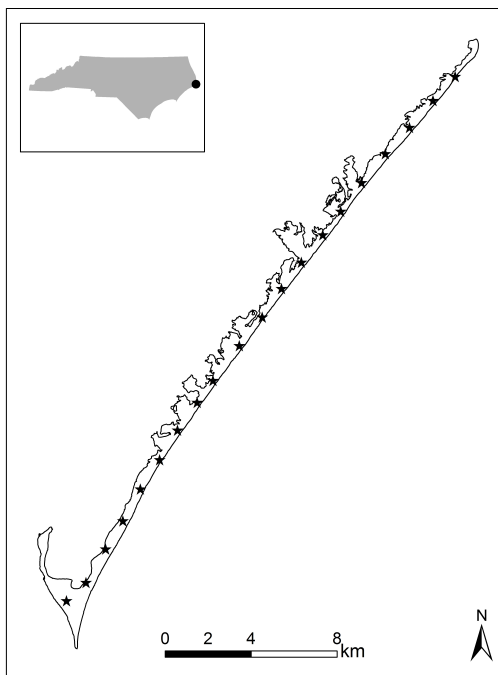


Figure 1.7: Camera traps (stars) set up on South Core Banks, a barrier island within Cape Lookout National Seashore, North Carolina (inset map) to estimate the raccoon population (see Chapt. ?? for details).

the state-space, these features need to be removed. You can create a precise state-space polygon in **ArcGIS** and read it into **R**, or create the polygon directly within **R**, by intersecting two shapefiles – one of the rectangle defining the outer limits of your state-space state and one of the landscape feature you want to remove. While you will most likely have to obtain the shapefile describing the landscape of and around your trapping grid (coastlines, water bodies etc.) from some external source, the polygon shapefile buffering your outermost trapping grid coordinates can easily be written in **R**.

If `xmin`, `xmax`, `ymin` and `ymax` mark the most extreme x and y coordinates of your trapping grid and b is the distance you want to buffer with, load the package **shapefiles** (Stabler, 2006) and issue the following **R** commands:

```

925 > xl= xmin-b
926 > xu= xmax+b
927 > yl= ymin-b
928 > yu= ymax+b
929
930 #create data frame with coordinate pairs
931 > dd <- data.frame(Id=c(1,1,1,1,1),X=c(xl,xu,xu,xl,xl),
932
```



```

943 + Y=c(y1,y1,yu,yu,y1))
944 > ddTable <- data.frame(Id=c(1),Name=c("Item1"))
945           #convert to shapefile, type polygon
946 > ddShapefile <- convert.to.shapefile(dd, ddTable, "Id", 5)
947           # name and save to location of choice
948 > write.shapefile(ddShapefile, 'c:/Test', arcgis=T)

```

949 You can read shapefiles into **R** loading the package `maptools` (Lewin-Koh
 950 et al., 2011) and using the function `readShapeSpatial()`. Make sure you read
 951 in shapefiles in UTM format, so that units of the trap array, the movement
 952 parameter σ and the state-space are all identical. Intersection of polygons can
 953 be done in **R** also, using the package `rgeos` (Bivand and Rundel, 2011) and
 954 the function `gIntersect()`. The area of your (single) polygon can be extracted
 955 directly from the state-space object `SSp`:

```

956 > area <- SSp@polygons[[1]]@Polygons[[1]]@area /1000000

```

957 Note that dividing by 1000000 will return the area in km^2 if your coordi-
 958 nates describing the polygon are in UTM. If your state-space consists of several
 959 disjunct polygons, you will have to sum the areas of all polygons to obtain the
 960 size of the state-space. To include this polygon into our MCMC sampler we
 961 need one last spatial **R** package, `sp` (Pebesma and Bivand, 2011), which has a
 962 function, `over()`, which allows us to check if a pair of coordinates falls within a
 963 polygon or not.⁴ All we have to do is embed this new check into the updating
 964 steps for the s_i :

```

965           #draw candidate value
966 Scand <- as.matrix(cbind(rnorm(M, S[,1], 2), rnorm(M, S[,2], 2)))
967           #convert to spatial points on UTM (m) scale
968 Scoord<-SpatialPoints(Scand*1000)
969           # check if scand is within the polygon
970 SinPoly<-over(Scoord,SSp)
971
972 for(i in 1:M) {
973   #if scand falls within polygon, continue update
974   if(is.na(SinPoly[i])==FALSE) {
975     ... [rest of the updating step remains the same]

```

976 Note that it is much more time-efficient to draw all M candidate values for
 977 s and check once if they fall within the state-space, rather than running the
 978 `over()` command for every individual pair of coordinates. To make sure that
 979 our initial values for s also fall within the polygon of \mathcal{S} , we use the function
 980 `runifpoint()` from the package `spatstat` (Baddeley and Turner, 2005), which

⁴Remember from Chapt. (??) that the `over` function takes as its second argument (among others) an object of the class "SpatialPolygons" or "SpatialPolygonsDataFrame". The former produces a vector while the latter produces a data frame (e.g., in the example above), which is important for how you index the output.

generates random uniform points within a specified polygon. You'll find this modified MCMC algorithm (`SCR0poisSSp`) in the **R** package `scrbook`.

Finally, observe that we are converting candidate coordinates of \mathcal{S} back to meters to match the UTM polygon. In all previous examples, for both the trap locations and the activity centers we have used UTM coordinates divided by 1000 to estimate σ on a km scale. This is adequate for wide ranging species like bears. In other cases you may center all coordinates on 0. No matter what kind of transformation you use on your coordinates, make sure to always convert candidate values for \mathcal{S} back to the original scale (UTM) before running the `over()` command.

1.8 Increasing Computational Speed

Using custom written MCMC algorithms in **R** is not only more flexible but can also be faster than using programs such as **JAGS** and especially **WinBUGS**. Also, **R** tends to use much less memory than **JAGS**, which can be crucial if you are running a large model but only have limited memory available. **WinBUGS** is limited in the amount of memory it can access and thus will likely not max out your memory, but as a trade-off, it will take a long time to run such models. In this chapter we have provided you with the guidelines to write your own MCMC sampler. But beyond the material that we have covered there are a number of ways you can make your sampler more efficient, through parallel computing or by accessing an alternative computer language such as **C++**. Exploring these options exhaustively is beyond the scope of this book; instead, in this section we will give you some pointers to get started with these more advanced computational issues.

1.8.1 Parallel computing

If you are using a computer with several cores, you can make use of parallel computing to speed up overall computation. In parallel computing we execute commands simultaneously on different cores of the computer, instead of running them serially on one single core. For example, imagine you have 4 cores available and you want to implement a for-loop in **R**; instead of going through the loop iteration by iteration, you can prompt **R** to execute iterations 1 to 4 simultaneously on the 4 different cores. The core that finishes first will then continue with iteration 5, and so on. There are several packages in **R** that allow you to induce parallel computing, such as `snow` (Tierney et al., 2011) and `snowfall` (Knaus, 2010), and the more current versions of **R** (from 2.14.0 upwards) come with a pre-installed set of functions grouped under the name `parallel`.

The MCMC algorithms developed here and in other parts of this book come with plenty of opportunities to parallelize computation. In various instances within the algorithm, we have for-loops across our augmented data set of size M , or we may have for-loops across sampling occasions. We also have for-loops across iterations of the algorithm, but since one iteration of the Markov

chain depends on the preceding iteration these should always be run serially, not in parallel. There is another dimension we can think of, and that is running multiple chains of an algorithm to assess convergence. This is a comparatively easy implementation of parallel computing and thus provides a good starting point to understand how it works in **R**.

Let's go back to the Ft. Drum black bear data we analyzed above with the cloglog version of the binomial SCR model (sec. 1.6) and run 3 parallel chains using `snowfall`. All we need to do is wrap our function `SCR0binom.cl` within another function that can then be executed in parallel, returning a list with one output matrix for each chain (install `snowfall` before executing the code below; we assume the data objects are already in your workspace from the previous analysis):

```
1034 > library(snowfall)
1035 ## create wrapper function
1036 > wrapper<-function(a){
1037 +   out<-SCR0binom.cl(y=Xaug, X=trapmat, M=M, x1=x1, xu=xu, y1=y1,
1038 +                     yu=yu, K=8, delta=c(0.1, 0.05, 2), niter=5000)
1039 +   return(out)
1040 + }
```

After creating the wrapper function we need to initialize the cluster of cores, defining that we want computation to be implemented in parallel and how many cores we want it to be run on. Here, we assume we have (at least) 3 cores, but if your computer only has 2, make sure to adjust the code accordingly (i.e., set `cpus=2`). In that case, 2 of the 3 chains will be run in parallel and whichever core finishes first will then pick up the third chain. Further, we have to export all **R** libraries and data to all the cores, and set up a random number generator, so that we do not get identical results from the different cores:

```
1049 > sfInit( parallel=TRUE, cpus=3 ) #initialize cluster
1050 > sfLibrary(scrbook) #export library scrbook
1051 > sfExportAll() #export all data in current workspace
1052 > sfClusterSetupRNG() #set up random number generator
1053 > outL=sfLapply(1:3,wrapper) # execute 'wrapper' 3 times
```

The object `outL` is a list of length 3, with one `out` matrix from the function `SCR0binom.cl` for each chain. After computation is complete, terminate the cluster using the command `sfStop()`. Note that the intermediate output of current values and acceptance rates in the **R** console is suppressed when using parallel computing. We can now look at the output as described previously using the package `coda`, by first defining `outL` to be a list of `mcmc` objects.

```
1060 > library(coda)
1061 #turn output into MCMC list
1062 > res<-mcmc.list(as.mcmc(outL[[1]]),as.mcmc(outL[[2]]),as.mcmc(outL[[3]]))
1063 > summary(window(res, start=1001)) #remove first 1000 iterations as burn-in
1064
1065 [... some output removed ...]
```

```

1066
1067           Mean          SD Naive SE Time-series SE
1068 sigma    1.9723  0.13093 0.0011952      0.0087055
1069 lam0     0.1115  0.01535 0.0001401      0.0009003
1070 psi      0.7130  0.10787 0.0009847      0.0077910
1071 N       499.6166 74.74934 0.6823650      5.4232653
1072
1073 2. Quantiles for each variable:
1074
1075           2.5%      25%      50%      75%      97.5%
1076 sigma    1.74339   1.8811   1.9637   2.0530   2.2618
1077 lam0     0.08443   0.1007   0.1105   0.1211   0.1438
1078 psi      0.52046   0.6350   0.7093   0.7814   0.9627
1079 N       366.00000 446.0000 497.0000 547.0000 674.0000

```

Now that we have parallel chains we can also use the function `gelman.diag` to evaluate if chains have converged:

```

1082 > gelman.diag(window(res, start=1001)) #assess chain convergence
1083

```

Potential scale reduction factors:

```

1085
1086           Point est. Upper C.I.
1087 sigma           1.01      1.04
1088 lam0            1.01      1.02
1089 psi             1.07      1.21
1090 N               1.07      1.21

```

Multivariate psrf

```

1093
1094 1.05

```

We can see that estimates are similar to what we observed when running a single chain (see sec. 1.6) and that all 3 chains appear to have converged, based on their point estimates of the \hat{R} statistic, but, as already noted before, for a real analysis we might want to run this model for quite a bit longer, to bring down the upper confidence interval limits on \hat{R} for ψ and N . If you have 3 cores then running these 3 parallel chains should not have taken longer than running a single chain. Yet if you look at the effective sample size now using `effectiveSize`, you can see that it has roughly tripled, as we would expect:

```

1103 > effectiveSize(window(res, start=1001))
1104
1105           sigma      lam0      psi      N
1106 272.6935 411.8384 167.4192 168.3355

```

1.8.2 Using C++

Parallel computing is a great tool to speed up computations, but its usefulness is limited by how many cores you have available. Even with a decent number

of cores, large models may still take a long time to run. A major reason for this is that for-loops in **R** are time consuming, whereas they are handled much more time efficiently in other computer languages such as **C++**. As we saw above, MCMC algorithms consist of for-loops within for-loops, so that it stands to reason that implementing them in a language like **C++** should make those algorithms run much faster. Being avid **R** users, we cannot claim to be fluent in **C++** or to be aware of all the opportunities this language brings for faster computing. It is also beyond the scope of this book to go into the nuts and bolts of how **C++** works or provide a tutorial, and we refer you to the vast amounts of online and print material designed to give the interested user an introduction to **C++**. Just google “introduction C++” and you are sure to come across sites such as <http://www.cplusplus.com> that provide step by step instructions to get you started. Here, we only want to point out one approach to linking **R** with **C++**: the packages **inline** (Sklyar et al., 2010) and **RcppArmadillo** (François et al., 2011). These two packages provide a very convenient interface between the two languages, but there are other ways of calling **C++** functions from within **R**, such as the `.Call` command. If you are interested, we suggest you refer to the package manuals and vignettes, as well as the online document “Writing R extensions” (at <http://cran.r-project.org/doc/manuals/R-exts.html>) for a much more thorough treatment of this topic.

In order to use **C++** you need a compiler such as **g++** that (together with other compilers, for example for **C** and **FORTRAN**) comes with **Rtools**, which you can easily download from the web (at <http://cran.r-project.org/bin/windows/Rtools/>). All of these compilers are part of the GNU compiler collection (<http://gcc.gnu.org/>). Make sure the version of **Rtools** matches your version of **R** or you may run into compilation errors later on. To give you a taste of **C++** we will show you how to write a function that calculates the squared distances of individual activity centers to all traps, as is implemented in the **scrbook** package in the function `e2dist` (to be exact, `e2dist` calculates the distance, not the squared distance), and compare performance between **R** and **C++**. We will refer to these functions as “distance functions”. First, let us set up dummy data – a matrix holding the coordinates of the trap array, outer limits of the state-space and uniformly distributed activity centers for $M = 700$ individuals:

```

> gx<-seq(1,10,1)
> gy<-seq(1,10,1)
> X<-as.matrix(expand.grid(gx, gy))
> M<-700
> J<-dim(X)[1]
> b<-3
> xl<-min(gx)-b
> xu<-max(gx)+b
> yl<-min(gy)-b
> yu<-max(gy)+b
> S<-cbind(runif(M, xl, xu), runif(M, yl,yu))

```

Next, we can write a “pedestrian” version of `e2dist` and check how long it

1157 takes to calculate the squared distance matrix:

```

1158 > Dfun<-function(M, J, S, X){
1159 + D2<-matrix(0, nrow=M, ncol=J)
1160 + for (i in 1:M){
1161 + for(j in 1:J){
1162 + D2[i,j]<-(S[i,1]-X[j,1])^2 + (S[i,2]-X[j,2])^2
1163 + }}
1164 + return(D2)
1165 + }
1166
1167 > system.time(
1168 + (D2R<-Dfun(M, J, S, X))
1169 + )
1170
1171     user  system elapsed
1172     0.81    0.01    0.82

```

1173 The code to implement the same function in **C++** using the `inline` and
1174 `RcppArmadillo` packages is shown in panel 1.3. These packages allow you to
1175 use a range of data formats such as lists and matrices, and they take care of
1176 compiling the code in **C++** and loading the resulting function into **R**. This is
1177 also referred to compiling **C++** code “on the fly”. You will see that the way
1178 the code is set up is reasonably similar to **R**. One difference that is worthy to
1179 point out is that in **C++** indices for vectors range from 0 to $n - 1$, NOT from
1180 1 to n , as in **R**. Note that with `inline` we only need to write the core of the
1181 code and define the type of the variables we want to pass to the function, while
1182 the `cxxfunction` call takes care of the rest. Once your function is compiled and
1183 loaded you should check out the full **C++** code by calling `DfunArma@code`.

1184 Executing this code shows that it is faster than the **R** version of the dis-
1185 tance function or `e2dist`; in fact it is too fast for the time resolution of the
1186 `system.time()` function to even give us a time estimate:

```

1187 > system.time(
1188 + (out<-DfunArma(M,J,S,X)))
1189
1190     user  system elapsed
1191     0      0      0

```

1192 While speed differences of less than 1 second may seem negligible, remem-
1193 ber that each command has to be executed at each iteration of the Markov
1194 chain. Especially with time-consuming models such as those for open popula-
1195 tions (Chapt. ??) or multi-session models (Chapt. ??) we believe that **C++**
1196 holds large potential to make implementation of such models more feasible.

1197 1.9 Summary and Outlook

1198 In a nutshell, programs like **JAGS** and **WinBUGS** do all the MCMC-related
1199 things that we went through in this chapter (and quite a bit more). Looking

```

### calculate squared distances using RcppArmadillo
library(inline)
library(RcppArmadillo)

#write core of function code
code<-'
/*define input, assign correct class (matrix, vector etc)*/
arma::mat Sn=Rcpp::as<arma::mat>(S);
arma::mat Xn=Rcpp::as<arma::mat>(X);
int Ntot=Rcpp::as<int>(M);
int ntraps=Rcpp::as<int>(J);
/*create matrix to hold squared distances*/
arma::mat D2(Ntot, ntraps);

/*loop over M and J to calculate distances*/
for (int i=0; i<Ntot; i++){
  for(int j=0; j<ntraps; j++){
    D2(i,j)= pow(Sn(i,0)-Xn(j,0), 2) + pow(Sn(i,1)-Xn(j,1), 2);
  }
}
/*return D2 in R format*/
return Rcpp::wrap(D2);
,

# compile and load
DfunArma<-cxxfunction(signature(M="integer", J="integer", S="numeric",
X="numeric"), plugin="RcppArmadillo", body=code)

```

Panel 1.3: Code to compute squared distance between individual activity centers and traps in **C++** from within **R** using **inline** and **RcppArmadillo**

1200 through your model, they determine which parameters they can use standard
 1201 Gibbs sampling for (i.e. for conjugate full conditional distributions). Then,
 1202 they determine whether to use adaptive rejection sampling, slice sampling or –
 1203 in the ‘worst’ case – Metropolis-Hastings sampling for the other full conditionals
 1204 (how the sampler is chosen differs among softwares). For MH sampling, they
 1205 will automatically tune the updater so that it works efficiently.

1206 Although these programs are flexible and extremely useful to perform MCMC
 1207 simulations, it sometimes is more efficient to develop your own MCMC algo-
 1208 rithm. Building an MCMC code follows three basic steps: Identify your model
 1209 including priors and express full conditional distributions for each model pa-
 1210 rameter. If full conditionals are parametric distributions, use Gibbs sampling
 1211 to draw candidate parameter values from those distributions; otherwise use

1212 Metropolis-Hastings sampling to draw candidate values from a proposal distri-
 1213 bution and accept or reject them based on their posterior probability densities.

1214 These custom-made MCMC algorithms give you more modeling flexibility
 1215 than existing software packages, especially when it comes to handling the state-
 1216 space: In **WinBUGS** and **JAGS** we define a continuous rectangular state-
 1217 space using the corner coordinates to constrain the uniform priors on the activity
 1218 centers \mathbf{s} . But what if a continuous rectangle is an inadequate description of the
 1219 state-space? In this chapter we saw that in **R** it only takes a few lines of code to
 1220 use any arbitrary polygon shapefile as the state-space, which is especially useful
 1221 when you are dealing with coastlines or large bodies of water that need remov-
 1222 ing from the state-space. Another example is the SCR **R** package **SPACECAP**
 1223 (Gopalaswamy et al., 2012) that was developed because implementation of an
 1224 SCR model with a discrete state-space was inefficient in **WinBUGS**.

1225 Another situations in which using a **BUGS** engine becomes increasingly
 1226 complicated or inefficient is when using point processes other than the homo-
 1227 geneous binomial point process (“uniformity of density”) which underlies the
 1228 basic SCR model (see sec. ?? in Chapt. ??). In Chapt. ?? you already saw an
 1229 example of an inhomogeneous point process model and we briefly introduce a
 1230 different point processes, implemented using a custom-made MCMC algorithm,
 1231 in Chapt. ?. Finally, Chapt. ? deals with partially marked populations us-
 1232 ing hand-made MCMC algorithms to handle the (partially) latent individual
 1233 encounter histories. While some of these models can be written in the **BUGS**
 1234 language, they are painstakingly slow; others (for example the classes of models
 1235 considered in Chapt. ?) cannot be implemented in **WinBUGS/JAGS** at all
 1236 and we have to either use likelihood based inference or develop our own MCMC
 1237 algorithms. In conclusion, while you can certainly get by using **BUGS/JAGS**
 1238 for standard SCR models, knowing how to write your own MCMC sampler gives
 1239 you more flexibility to tailor these models to your specific needs.

Bibliography

- Baddeley, A. and Turner, R. (2005), “Spatstat: an R package for analyzing spatial point patterns,” *Journal of Statistical Software*, 12, 1–42, ISSN 1548-7660.
- Bivand, R. and Rundel, C. (2011), *rgeos: Interface to Geometry Engine - Open Source (GEOS)*, r package version 0.1-8.
- Buckland, S., Anderson, D., Burnham, K., Laake, J., Borchers, D., and L, T. (2001), *Introduction to distance sampling: estimating abundance of biological populations*, Oxford, UK: Oxford University Press.
- Casella, G. and George, E. I. (1992), “Explaining the Gibbs sampler,” *American Statistician*, 46, 167–174.
- François, R., Eddelbuettel, D., and Bates, D. (2011), *RcppArmadillo: Rcpp integration for Armadillo templated linear algebra library*, r package version 0.2.25.
- Gelfand, A. and Smith, A. (1990), “Sampling-based approaches to calculating marginal densities,” *Journal of the American statistical association*, 85, 398–409.
- 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.
- Geman, S. and Geman, D. (1984), “Stochastic relaxation, Gibbs distributions, and the Bayesian restoration of images,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, PAMI-6, 721–741.
- Gilks, W. and Wild, P. (1992), “Adaptive rejection sampling for Gibbs sampling,” *Applied Statistics*, 41, 337–348.
- Gilks, W. R., Thomas, A., and Spiegelhalter, D. J. (1994), “A Language and Program for Complex Bayesian Modelling,” *Journal of the Royal Statistical Society. Series D (The Statistician)*, 43, 169–177, ArticleType: primary_article / Issue Title: Special Issue: Conference on Practical Bayesian Statistics, 1992 (3) / Full publication date: 1994 / Copyright 1994 Royal Statistical Society.

- 1271 Gopalaswamy, A. M., Royle, A. J., Hines, J., Singh, P., Jathanna, D., Ku-
1272 mar, N. S., and Karanth, K. U. (2012), “Program SPACECAP: software for
1273 estimating animal density using spatially explicit capturerecapture models,”
1274 *Methods in Ecology and Evolution*, online early, r package version 1.0.4.
- 1275 Hastings, W. (1970), “Monte Carlo sampling methods using Markov chains and
1276 their applications,” *Biometrika*, 57, 97–109.
- 1277 Knaus, J. (2010), *snowfall: Easier cluster computing (based on snow)*., r package
1278 version 1.84.
- 1279 Lewin-Koh, N. J., Bivand, R., contributions by Edzer J. Pebesma, Archer, E.,
1280 Baddeley, A., Bibiko, H.-J., Dray, S., Forrest, D., Friendly, M., Giraudoux, P.,
1281 Golicher, D., Rubio, V. G., Hausmann, P., Hufthammer, K. O., Jagger, T.,
1282 Luque, S. P., MacQueen, D., Niccolai, A., Short, T., Stabler, B., and Turner,
1283 R. (2011), *maptools: Tools for reading and handling spatial objects*, r package
1284 version 0.8-10.
- 1285 Link, W. A. and Barker, R. J. (2010), *Bayesian Inference: With Ecological*
1286 *Applications*, London, UK: Academic Press.
- 1287 Metropolis, N., Rosenbluth, A., Rosenbluth, M., Teller, A., Teller, E., et al.
1288 (1953), “Equation of state calculations by fast computing machines,” *The*
1289 *journal of chemical physics*, 21, 1087–1092.
- 1290 Metropolis, N. and Ulam, S. (1949), “The Monte Carlo method,” *Journal of the*
1291 *American Statistical Association*, 44, 335–341.
- 1292 Neal, R. (2003), “Slice sampling,” *Annals of Statistics*, 31, 705–741.
- 1293 Pebesma, E. and Bivand, R. (2011), *Package ‘sp’*, r package version 0.9-91.
- 1294 Plummer, M. (2003), “JAGS: A program for analysis of Bayesian graphical mod-
1295 els using Gibbs sampling,” in *Proceedings of the 3rd International Workshop*
1296 *on Distributed Statistical Computing (DSC 2003)*. March, pp. 20–22.
- 1297 Plummer, M., Best, N., Cowles, K., and Vines, K. (2006), “CODA: Convergence
1298 Diagnosis and Output Analysis for MCMC,” *R News*, 6, 7–11.
- 1299 Robert, C. P. and Casella, G. (2004), *Monte Carlo statistical methods*, New
1300 York, USA: Springer.
- 1301 — (2010), *Introducing Monte Carlo Methods with R*, New York, USA: Springer.
- 1302 Roberts, G. O. and Rosenthal, J. S. (1998), “Optimal scaling of discrete ap-
1303 proximations to Langevin diffusions,” *Journal of the Royal Statistical Society:*
1304 *Series B (Statistical Methodology)*, 60, 255–268.
- 1305 Sklyar, O., Murdoch, D., Smith, M., Eddelbuettel, D., and François, R. (2010),
1306 *inline: Inline C, C++, Fortran function calls from R*, r package version 0.3.8.

- 1307 Sollmann, R., Gardner, B., Parsons, A., Stocking, J., McClintock, B., Simons,
1308 T., Pollock, K., and O'Connell, A. (2013), "A spatial mark-resight model
1309 augmented with telemetry data," *Ecology*.
- 1310 Stabler, B. (2006), *shapefiles: Read and Write ESRI Shapefiles*, r package version
1311 0.6.
- 1312 Tierney, L., Rossini, A. J., Li, N., and Sevcikova, H. (2011), *snow: Simple*
1313 *Network of Workstations*, r package version 0.3-7.