

# 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.

### 1.0.1 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** provides 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 shape file in a **BUGS** model. Similarly, models of space usage that take into account ecological distance (Chapt. ??) cannot be implemented in the **BUGS** engines. Moreover, there are classes of SCR models that we have not been able to implement effectively using likelihood methods, and are inefficient to run in the **BUGS** engines. Examples of those models are covered in Chaps. ?? and ??.

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.

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. Not to mention that, if you are an **R** junkie, it can actually be fun. 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.

## 1.1 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  $\theta$ , 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  $\theta$  is that it can be used to make probability statements about  $\theta$ , such as the probability that  $\theta$  is equal to some value, or the probability that  $\theta$  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

70 to directly compute the posterior. Recall Bayes' theorem:

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

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

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

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

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

78 where  $\sigma$  is known and our objective is to obtain an estimate of  $\mu$ . To fully specify  
 79 the model in a Bayesian framework, we first have to define a prior distribution  
 80 for  $\mu$ . Recall from Chapt. ?? that for certain data models, certain priors lead  
 81 to conjugacy, i.e. if you choose a certain prior for your parameter, the posterior  
 82 distribution will be of a known parametric form. The conjugate prior for the  
 83 mean of a Normal model is also a Normal distribution:

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

84 If  $\mu_0$  and  $\sigma_0^2$  are fixed, the posterior for  $\mu$  has the following form (for some of  
 85 the algebra behind this, see Chapt. 2 in Gelman et al. (2004)):

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

86 where

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

87 And

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

88 We can directly obtain estimates of interest from this Normal posterior distri-  
 89 bution, such as the mean  $\hat{\mu}$  and its variance; we do not need to apply MCMC,  
 90 since we can recognize the posterior as a parametric distribution, including the  
 91 normalizing constant  $[y]$ . But generally we will be interested in more complex  
 92 models with several, say  $m$ , parameters. In this case, computing  $[y]$  from Eq.  
 93 1.1 requires  $m$ -dimensional integration, which can be difficult or impossible.  
 94 Thus, the posterior distribution is generally only known up to a constant of  
 95 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  $\theta$  (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 chain<sup>1</sup>. 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  $\theta$ , its standard error and confidence bounds, can be obtained directly from this approximation of the posterior.

## 1.2 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.

### 1.2.1 Gibbs sampling

Gibbs sampling was named after the physicist J.W. Gibbs by Geman and Geman (1984), who applied the algorithm to a Gibbs distribution<sup>2</sup>. 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  $\mu$ , the posterior distribution of  $\mu|y$  is also Normal. Conversely, with a fixed (known)  $\mu$ , but unknown variance, the

<sup>1</sup>Remember 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,  $\theta^{(t-1)}$ .

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

conjugate prior for  $\sigma^2$  is an Inverse-Gamma distribution with shape and scale parameters  $a$  and  $b$ :

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

With fixed  $a$  and  $b$ , the posterior  $[\sigma^2|\mu, y]$  is also an Inverse-Gamma distribution, namely:

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

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 neither  $\mu$  nor  $\sigma^2$ , which is probably the more common case? The joint posterior distribution of  $\mu$  and  $\sigma^2$  now has the general structure

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

or

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

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

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

The full conditional of  $\mu$  again takes the form of the Normal distribution shown in Eq. 1.2; similarly,  $[\sigma^2|\mu, y]$  takes the form of the Inverse-Gamma distribution shown in Eq. 1.3, both distribution we can easily sample from. And this is precisely what we do when using Gibbs sampling: we break down high-dimensional problems into convenient one-dimensional problems by constructing the full conditional distributions for each model parameter separately; and we sample from these full conditionals, which, if we choose conjugate priors, are known parametric distributions. Let's put the concept of Gibbs sampling into the MCMC framework of generating successive samples, using our simple Normal model with unknown  $\mu$  and  $\sigma^2$  and conjugate priors as an example. These are the steps you need in order to build a Gibbs sampler:

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

**Step 1:** Draw  $\theta_1^{(1)}$  from the conditional distribution  $[\theta_1^{(1)}|\theta_2^{(0)}, \dots, \theta_d^{(0)}]$ . Here,  $\theta_1$  is  $\mu$ , which we draw from the Normal distribution in Eq. 1.2 using  $\sigma^{(0)}$  as value for  $\sigma$ .

**Step 2:** Draw  $\theta_2^{(1)}$  from the conditional distribution  $[\theta_2^{(1)}|\theta_1^{(1)}, \theta_3^{(0)}, \dots, \theta_d^{(0)}]$ . Here,  $\theta_2$  is  $\sigma$ , which we draw from the Inverse-Gamma distribution of Eq. 1.3, using  $\mu^{(1)}$  as value for  $\mu$ .

167 **Step 3, ..., d:** Draw  $\theta_3^{(1)}, \theta_4^{(1)}, \dots, \theta_d^{(1)}$  from their conditional distribution  
 168  $[\theta_3^{(1)} | \theta_1^{(1)}, \theta_2^{(1)}, \theta_4^{(0)}, \dots, \theta_d^{(0)}], \dots, [\theta_d^{(1)} | \theta_1^{(1)}, \dots, \theta_{d-1}^{(1)}]$ . In our example we have  
 169 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 In terms of **R** coding, this means we have to write Gibbs updaters for  $\mu$  and  
 172  $\sigma^2$  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 according to Eq. 7.2
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 according to Eq. 7.3
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  $\mu$  and  $\sigma$  and conjugate priors (Normal and Inverse-Gamma, respectively) for both parameters.

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

```

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

```

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

```

195 plot(out[,1], out[,2])

```

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

```

199 par(mfrow=c(1,2))
200 hist(out[,1]); hist(out[,2])

```

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

```

216 summary(mod[501:10000,])
217      mu      sig
218 Min.   :4.935  Min.   :0.4652
219 1st Qu.:4.988  1st Qu.:0.4930
220 Median :4.998  Median :0.5006

```

```

221 Mean      :4.998      Mean      :0.5008
222 3rd Qu.   :5.009      3rd Qu.   :0.5084
223 Max.     :5.062      Max.     :0.5486

```

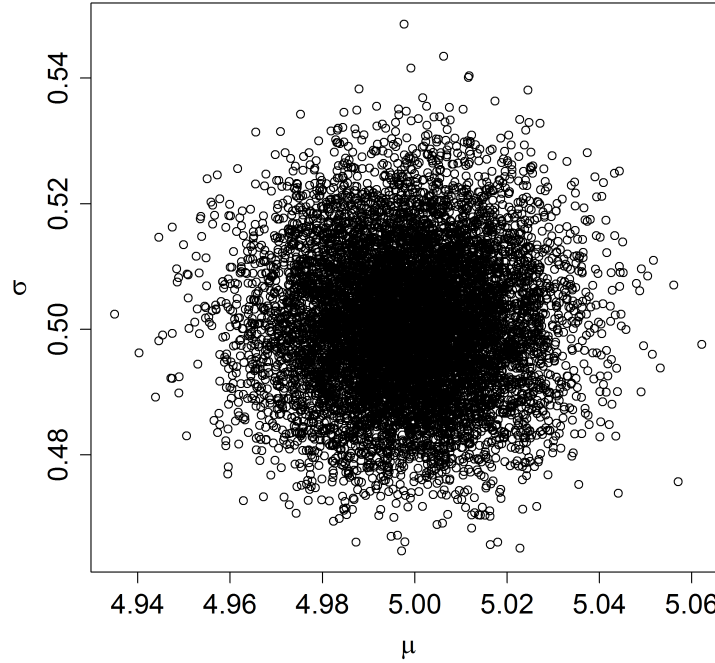


Figure 1.1: Joint posterior distribution of  $\mu$  and  $\sigma$  from a Normal Model

### 1.2.2 Metropolis-Hastings sampling

Although it is applicable to a wide range of problems, the limitations of Gibbs sampling are immediately obvious: what if we do not want to use conjugate priors or what if we cannot recognize the full conditional distribution as a parametric distribution, or simply do not want to worry about these issues? The most general solution is to use the Metropolis-Hastings (MH) algorithm, which also goes back to the work by Metropolis et al. (1953). You saw the basics of this algorithm in Chapt. ???. In a nutshell, because we do not recognize the posterior  $[\theta|y]$  as a parametric distribution, the MH algorithm generates samples from a known proposal distribution, say  $h(\theta)$ , that depends on the value of  $\theta$  at the previous time step,  $\theta^{t-1}$ . The candidate value  $\theta^*$  is accepted with probability.



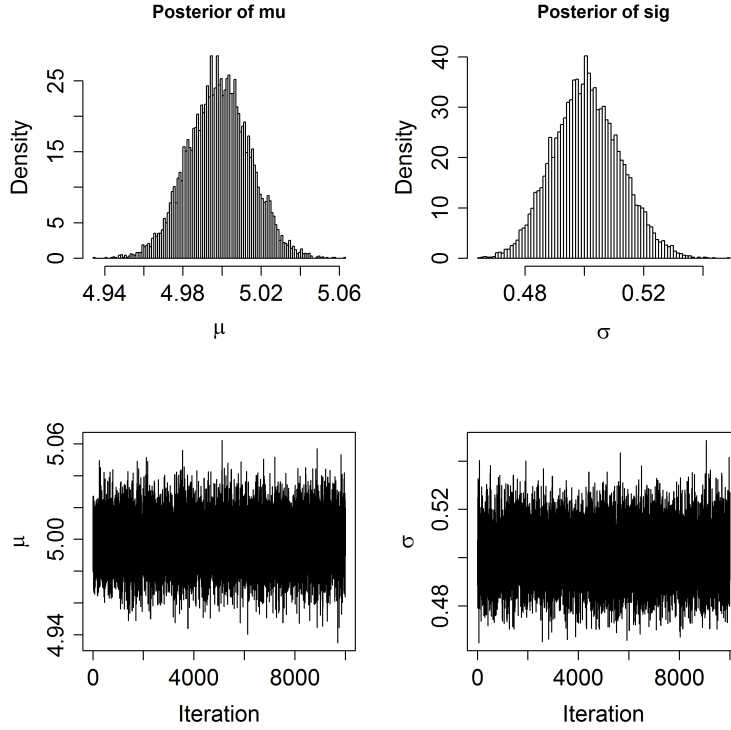


Figure 1.2: Plots of the posterior distributions of  $\mu$  (upper left) and  $\sigma$  (upper right) from a Normal model and time series plots of  $\mu$  (lower left) and  $\sigma$  (lower right).

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

236 Proposal distributions can be absolutely anything! You can generate candi-  
 237 date values from a Normal(0,1) distribution, from a Uniform(-3455,3455) distri-  
 238 bution, or anything of proper support. Note, however, that good choices of  $h(\cdot)$   
 239 are those that approximate the posterior distribution. Obviously if  $h(\cdot) = [\theta|y]$   
 240 (i.e., the posterior) then you always accept the draw, and it stands to reason  
 241 that proposals that are more similar to  $[\theta|y]$  will lead to higher acceptance  
 242 probabilities.

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

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

244 In that case these two terms just cancel out from the MH acceptance probability  
 245 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  $\delta$ , 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,  $\theta^*$  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

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

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  $\theta$  from  $h(\theta^*|\theta^{t-1})$ . We often use a Normal proposal distribution, so we draw  $\theta^{(1)}$  from  $\text{Normal}(\theta^{(0)}, \delta)$ , where  $\delta$  is the variance 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  $\theta$ :

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

In our example,

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

281 **Step 3:** Set

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

282 We can do this last step by drawing a random number  $u$  from a Uniform(0, 1)  
283 and accept  $\theta^*$  if  $u < r$ . Repeat for  $t = 1, 2, \dots$  a large number of samples. The  
284 **R** code for this MH sampler is provided in Panel 1.2.

---

```
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.

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

291 Comparing MH sampling to Gibbs sampling, where all draws from the con-  
292 ditional distribution are used, in the MH algorithm we discard a portion of the  
293 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,  $\delta$ , 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.  $\delta$  should be chosen (a) large enough so that each step of drawing a new proposal value for  $\theta$  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 automatically adjusted until it reaches a user-defined acceptance rate, at which point the adaptive phase ends and the actual Markov chain begins. This is computationally a little more advanced. Link and Barker (2010) discuss this in more detail. It is important the samples drawn during the adaptive phase are discarded. To illustrate the effects of tuning, we ran the Metropolis-within-Gibbs algorithm in Panel 1.2 with  $\delta = 0.01$ ,  $\delta = 0.2$  and  $\delta = 1$ . The first 150 iterations for  $\theta$  are shown in Fig. 1.3. We see that for a very small  $\delta$  (the dashed line) the burn-in is extremely slow - after 150 iterations the chain isn't even half way there, while for the other two values of  $\delta$  (solid and dotted) the burn-in phase seems to be over after only about 10 iterations. While  $\delta = 0.2$  leads to reasonably good mixing, the chain clearly gets stuck on certain values with  $\delta = 1$ .

Other than graphically, you can easily check acceptance rates for the parameters you monitor (that are part of your output) using the `rejectionRate()` function of the package `coda` (we will talk more about this package a little later on). Do not let the term 'rejection rate' confuse you; it is simply  $1 - \text{acceptance rate}$ . There may be parameters – for example, individual values of a random effect or latent variables – that you do not want to save, though, and in our next example we will show you a way to monitor their acceptance rates with a few extra lines of code.

### 1.2.3 Metropolis-within-Gibbs

One weakness of the MH sampler is that formulating the joint posterior when evaluating whether to accept or reject the candidate values for  $\theta$  becomes increasingly complex or inefficient as the number of parameters in a model increases. As you already saw in Chapt. ??, in these cases you can simply combine MH sampling and Gibbs sampling. You can use Gibbs sampling to break down your high-dimensional parameter space into easy-to-handle one-dimensional conditional distributions and use MH sampling for these conditional

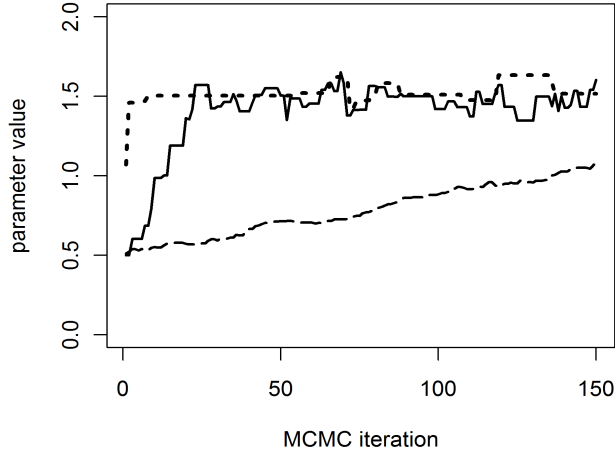


Figure 1.3: Time series plots of  $\theta$  from a MH algorithm with tuning parameter  $\delta = 0.01$  (dashed line),  $0.2$  (solid line) and  $1$  (dotted line).

distributions. Better yet, if you have some conjugacy in your model, you can use the more efficient Gibbs sampling for these parameters and one-dimensional MH for all the others. You have already seen the basics of how to build both types of algorithms, so we can jump straight into an example here and build a Metropolis-within-Gibbs algorithm.

**GLMMs: Poisson regression with a random effect** Let's assume a model that gets us closer to the problem we ultimately want to deal with - a GLMM. 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 use Normal priors on  $\alpha$  and  $\beta$ ,

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

and

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

In this model, we do not specify  $\mu_\alpha$  and  $\sigma_\alpha$ , 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  $\beta$  is constant across the  $i$  study sites, and for analysis we would set  $\mu_\beta$  and  $\sigma_\beta$ . 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  $\alpha_i$  are constructed similarly:

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

and so on for all elements of  $\alpha$ . The full-conditional for  $\beta$  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  $\alpha$ :

$$\begin{aligned}[\mu_\alpha | \alpha] &\propto [\alpha | \mu_\alpha, \sigma_\alpha][\mu_\alpha] \\ [\sigma_\alpha | \alpha] &\propto [\alpha | \mu_\alpha, \sigma_\alpha][\sigma_\alpha]\end{aligned}$$

Since we assumed  $\alpha$  to come from a Normal distribution, the choice of priors for  $\mu_\alpha$  (Normal) and  $\sigma_\alpha^2$  (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  $\alpha$  and  $\beta$  we use Normal proposal distributions with standard deviations  $\delta_\alpha$  and  $\delta_\beta$ .

First, we set the initial values  $\alpha^{(0)}$  and  $\beta^{(0)}$ . Then, starting with  $\alpha_1$ , we draw  $\alpha_1^{(1)}$  from  $\text{Norm}(\alpha_1^{(0)}, \delta_\alpha)$ , 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  $\alpha$ .

For  $\beta$ , we draw  $\beta^{(1)}$  from  $\text{Norm}(\beta^{(0)}, \delta_\beta)$ , 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  $\mu_\alpha$  and  $\sigma_\alpha^2$ , we sample directly from the full conditional distributions (Eq. 1.2 and Eq. 1.3):

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

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  $\alpha$ , which we updated before, and  $n_\alpha$  is the length of  $\alpha$ . For  $\sigma_\alpha^2$  we use  $\sigma_\alpha^{2(1)} \sim \text{Inverse-Gamma}(a_n, b_n)$ , where

$$a_n = n_\alpha/2 + a_0, \text{ and } b_n = 0.5 \sum_{i=1}^{n_\alpha} (\alpha_i^{(1)} - \mu_\alpha^{(1)})^2 + b_0.$$

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

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

```
#initiate counter for acceptance rate of alpha
alphaUps<-0

#loop over sites, update intercepts alpha one at a time;
#only data at site i contributes information
#lev is the number of sites i
for (i in 1:lev) {
  alpha.cand<-rnorm(1, alpha[i], delta_alpha)
  loglike<- sum(dpois (y[site==i], exp(alpha[i] + beta*x[site==i]),
    log=TRUE))
  logprior<- dnorm(alpha[i], mu_alpha,sig_alpha, log=TRUE)
```

```

410 loglike.cand<- sum(dpois (y[site==i], exp(alpha.cand + beta *x[site==i]),
411   log=TRUE))
412 logprior.cand<- dnorm(alpha.cand, mu_alpha,sig_alpha, log=TRUE)
413 if (runif(1)< exp((loglike.cand+logprior.cand) -(loglike+logprior))) {
414   alpha[i]<-alpha.cand
415   alphaUps<-alphaUps+1
416 }
417 }
418
419 #lets you check the acceptance rate of alpha at every 100th iteration
420 if(iter %% 100 == 0) {
421   cat("   Acceptance rates\n")
422   cat("   alpha =", alphaUps/lev, "\n")
423 }

```

#### 424 1.2.4 Rejection sampling and slice sampling

425 While MH and Gibbs sampling are probably the most widely applied algorithms  
 426 for posterior approximation, there are other options that work under certain  
 427 circumstances and may be more efficient when applicable. **WinBUGS** applies  
 428 these algorithms and we want you to be aware that there is more out there  
 429 to approximate posterior distributions than Gibbs and MH. One alternative  
 430 algorithm is rejection sampling. Rejection sampling is not an MCMC method,  
 431 since each draw is independent of the others. The method can be used when  
 432 the posterior  $[\theta|y]$  is not a known parametric distribution but can be expressed  
 433 in closed form. Then, we can use a so-called envelope function, say,  $g(\theta)$ , that  
 434 we can easily sample from, with the restriction that  $[\theta|y] < M \times g(\theta)$ . We then  
 435 sample a candidate value for  $\theta$  from  $g(\theta)$ , calculate  $r = [\theta|y]/M \times g(\theta)$  and keep  
 436 the sample with the probability  $r$ .  $M$  is a constant that has to be picked so that  
 437  $r$  lies between 0 and 1, for example by evaluating both  $[\theta|y]$  and  $g(\theta)$  at  $n$  points  
 438 and looking at their ratios. Rejection sampling only works well if  $g(\theta)$  is similar  
 439 to  $[\theta|y]$ , and packages like **WinBUGS** use adaptive rejection sampling (Gilks  
 440 and Wild, 1992), where a complex algorithm is used to fit an adequate and  
 441 efficient  $g(\theta)$  based on the first few draws. Though efficient in some situations,  
 442 rejection sampling does not work well with high-dimensional problems, since  
 443 it becomes increasingly hard to define a reasonable envelope function. For an  
 444 example of rejection sampling in the context of SCR models, see Chapt. ??,  
 445 where we use it to simulation non-stationary point processes.

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

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

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

451 Slice sampling can be applied in many situations; however, implementing an  
 452 efficient slice sampling procedure can be complicated. We refer the interested



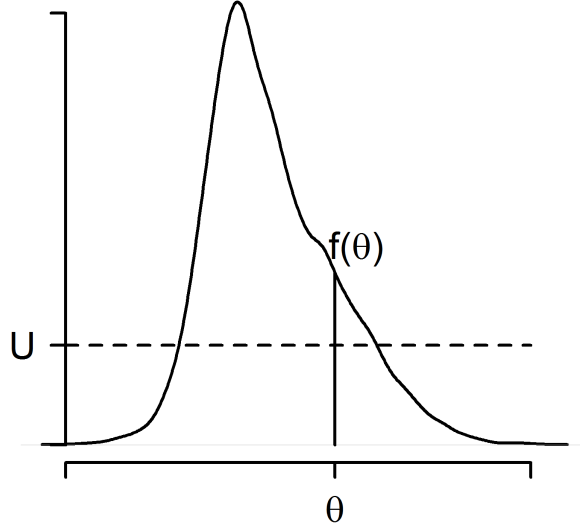


Figure 1.4: Slice sampling. For  $U \sim \text{Unif}(0, [\theta|y])$ , we can sample  $\theta$  from the vertical slice of  $[\theta|y]$  at  $U$ ;  $\theta|U \sim \text{Unif}(B)$ , where  $B = \{\theta : [\theta|y] \geq U\}$ .

453 reader to Robert and Casella (2010, Chapt. 7) for a simple example. Both rejection  
 454 sampling and slice sampling can be applied on one-dimensional conditional  
 455 distributions within a Gibbs sampling setup.

### 456 1.3 MCMC for closed capture-recapture Model 457 Mh

458 By now you have seen MCMC samplers for some simple GL(M)M's. Now, to  
 459 ease you into more complex models, we construct our own MCMC algorithm using  
 460 a Metropolis-within-Gibbs sampler for the non-spatial model with individual  
 461 heterogeneity in capture probability, model  $M_h$ , developed in Chapt. ??.

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

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

466 The binomial observation model is expressed conditional on the latent variables  
 467  $z_i$ .

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

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

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

As usual, we 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

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

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

For prior distributions, we assume that  $\mu_p, \sigma_p, \psi$  are mutually independent and for  $\mu_p$  and  $\sigma_p$  we use improper uniform priors, and  $\psi \sim \text{Unif}(0, 1)$ . 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  $\psi$ ,  $\mu_p$ , or  $\sigma_p$ . As such, the full-conditional for the structural parameter  $\psi$  only depend on the collection of data augmentation variables  $z_i$ , and that for  $\mu_p$  and  $\sigma_p$  will only depend on the collection of latent variables  $p_i; i = 1, 2, \dots, M$ . The full conditionals for all the unknowns are as follows:

(1) For  $p_i$ :

$$[p_i|y_i, \mu_p, \sigma_p, z_i] \propto \begin{cases} [y_i|p_i][p_i|\mu_p, \sigma_p^2] & \text{if } z_i = 1 \\ [p_i|\mu_p, \sigma_p] & \text{if } z_i = 0 \end{cases}$$

(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  $\mu_p$ :

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

(4) For  $\sigma_p$ :

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

(5) For  $\psi$ :

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

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  $x_i \sim \text{Bernoulli}(p)$  and  $p \sim \text{Beta}(a, b)$  is

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

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 term in the full conditionals for  $\mu_p$  and  $\sigma_p$  reflects the improper prior we chose for both parameters). 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.

```
### update the logit(p) parameters
lp.cand<- rnorm(M,lp,1) # 1 is a tuning parameter
p.cand<-plogis(lp.cand)
ll<-dbinom(ytot,K,z*p, log=T)
prior<-dnorm(lp,mu,sigma, log=T)
llcand<-dbinom(ytot,K,z*p.cand, log=T)
prior.cand<-dnorm(lp.cand,mu,sigma, log=T)

kp<- runif(M) < exp((llcand+prior.cand)-(ll+prior))
p[kp]<-p.cand[kp]
lp[kp]<-lp.cand[kp]
```

The parameters  $\mu_p$  and  $\sigma_p$  are also updated using MH steps (see the code for  $\mu_p$  below). In truth, we could also sample  $\mu_p$  and  $\sigma_p^2$  directly with certain choices of prior distributions. For example, if  $\mu_p \sim \text{Normal}(0,1000)$  then the full conditional for  $\mu_p$  is also Normal (see sec. 1.2.1), etc..

```
p0.cand<- rnorm(1,p0,.05)
if(p0.cand>0 & p0.cand<1){
mu.cand<-log(p0.cand/(1-p0.cand))
ll<-sum(dnorm(lp,mu,sigma,log=TRUE))
llcand<-sum(dnorm(lp,mu.cand,sigma,log=TRUE))
if(runif(1)<exp(llcand-ll)) {
mu<-mu.cand
p0<-p0.cand
}
}
```

For  $\psi$  we can easily sample directly from the beta distribution:

```
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 `z.cand<-ifelse(z==1,0,1)`. You can check out the full code by invoking `modelMh()` from the **R** package `scrbook`.

## 1.4 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'll 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  $\lambda_{ij}$ ,

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

The link between individual location, movement and trap encounter rates is made by the assumption that  $\lambda_{ij}$ , 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 half-normal form

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

where  $\lambda_0$  is the baseline trap encounter rate at  $d_{ij} = 0$  and  $\sigma$  controls the shape of the half-normal function.

In order to estimate the number of  $\mathbf{s}_i$  in  $\mathcal{S}$  (or any subset of  $\mathcal{S}$ ),  $N$ , we use data augmentation (sec. ??) and create  $M - n$  all-zero encounter histories, where  $n$  is the number of individuals we observed and  $M$  is a somewhat arbitrary number

that is larger than  $N$ . We estimate  $N$  by summing over the auxiliary data augmentation variables,  $z_i$ , which is 1 if the individual is part of the population and 0 if not, and assume that  $z_i$  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:

$\psi$ : the  $\text{Uniform}(0, 1)$  is required as part of the data augmentation procedure and in general is a natural choice of an uninformative prior for a probability. It will also lead to conjugacy as we saw in the example of model  $M_h$ , so that we can update  $\psi$  directly from its full conditional distribution using Gibbs sampling.

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

$\sigma$ : we can conceive several priors for  $\sigma$  but let's assume an improper prior, one that is Uniform over  $(-\infty, \infty)$ . We will see why this is convenient when we construct the full conditionals for  $\sigma$ .

$\lambda_0$ : analogous, we will use a  $\text{Uniform}(-\infty, \infty)$  improper prior for  $\lambda_0$ .

The parameter that is the objective of our modeling,  $N$ , is a derived parameter that we can obtain by summing all  $z_i$ :

$$N = \sum_{i=1}^M z_i$$

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

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

Since the improper prior implies that  $[\sigma] \propto 1$ , we can reduce this further to

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

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

```

590 sig.cand <- rnorm(1, sigma, 0.1) #draw candidate value
591 if(sig.cand>0){ #automatically reject sig.cand that are <0
592   lam.cand <- lam0*exp(-(d*d)/(2*sig.cand*sig.cand))
593   ll<- sum(dpois(y, lam*z, log=TRUE))
594   llcand <- sum(dpois(y, lam.cand*z, log=TRUE))
595   if(runif(1) < exp( llcand - ll) ){
596     ll<-llcand
597     lam<-lam.cand
598     sigma<-sig.cand
599   }
600 }

```

601 These steps are analogous for  $\lambda_0$  and  $\mathbf{s}_i$  and we will use MH steps for all of  
 602 these parameters. Similar to the random intercepts in our Poisson GLMM, we  
 603 update each  $\mathbf{s}_i$  individually. Note that to be fully correct, the full conditional for  
 604  $\mathbf{s}_i$  contains both the likelihood and prior component, since we did not specify  
 605 an improper, but a proper Uniform prior on  $\mathbf{s}_i$ . However, with a Uniform  
 606 distribution the probability density of any value is  $1/(\text{upper limit} - \text{lower limit})$   
 607 = constant. Thus, the prior components are identical for both the current  
 608 and the candidate value and can be ignored (formally, when you calculate the  
 609 ratio of posterior densities,  $r$ , the identical prior component appears both in the  
 610 numerator and denominator, so that they cancel each other out).

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

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

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

```

614 zUps <- 0 #set counter to monitor acceptance rate
615 for(i in 1:M) {
616   #no need to update seen individuals, since their z =1
617   if(seen[i])
618     next
619   zcand <- ifelse(z[i]==0, 1, 0)
620   llz <- sum(dpois(y[i,], lam[i,]*z[i], log=TRUE))
621   llcand <- sum(dpois(y[i,], lam[i,]*zcand, log=TRUE))
622
623   prior <- dbinom(z[i], 1, psi, log=TRUE)
624   prior.cand <- dbinom(zcand, 1, psi, log=TRUE)
625   if(runif(1) < exp((llcand+prior.cand)-(llz+prior))){
626     z[i] <- zcand
627     zUps <- zUps+1
628   }
629 }

```

630 The parameter  $\psi$  is a hyperparameter of the model, with an uninformative prior  
 631 distribution of Uniform(0, 1) or Beta(1, 1), so that

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

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

#### 635 1.4.1 SCR model with binomial encounter process

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

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

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

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

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

```

648     lam0.cand <- rnorm(1, lam0, 0.1)
649     #automatically reject lam0.cand that are <0
650     if(lam0.cand >0){
651         lam.cand <- lam0.cand*exp(-(d*d)/(2*sigma*sigma))
652         p.cand <- 1-exp(-lam.cand)
653         ll<- sum(dbinom(y, K, pmat *z, log=TRUE))
654         llcand <- sum(dbinom(y, K, p.cand *z, log=TRUE))
655         if(runif(1) < exp( llcand - ll) ){
656             ll<-llcand
657             pmat<-p.cand
658             lam0<- lam0.cand
659         }
660     }
```

661 Another possibility is to model variation in the individual and site specific  
 662 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:

```

668   p0.cand <- rnorm(1, p0, 0.1)
669   if(p0.cand > 0 & p0.cand < 1 ){
670     #automatically rejects lam0.cand that are not {0,1}
671     p.cand <- p0.cand*exp(-(d*d)/(2*sigma*sigma))
672     ll<- sum(dbinom(y, K, pmat *z, log=TRUE))
673     llcand <- sum(dbinom(y, K, p.cand *z, log=TRUE))
674     if(runif(1) < exp( llcand - ll) ){
675       ll<-llcand
676       pmat<-p.cand
677       p0<- p0.cand
678     }
679   }

```

#### 1.4.2 Looking at model output

Now that you have an MCMC algorithm to analyze spatial capture-recapture data with, let's run an actual analysis so we can look at the output. As an example, we will use the Fort Drum bear data set we first introduced in Chapt. ?? and already analyzed in Chapt. ?? with traditional non-spatial models (and that you will see again in Chapt. ??). You can load the Fort Drum data (`data(beardata)`), extract the trap locations (`trapmat`) and detection data (`bearArray`) and build the augmented  $M \times J$  array of individual encounter histories:

```

689 M=700
690 trapmat<-beardata$trapmat
691 #summarizes captures across occasions
692 bearmat<-apply(beardata$bearArray, 1:2, sum)
693 Xaug<-matrix(0, nrow=M, ncol=dim(trapmat)[1])
694 Xaug[1:dim(bearmat)[1],]<-bearmat #create augmented data set

```

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

```

698 x1<- min(trapmat[,1])- 20
699 y1<- min(trapmat[,2])- 20
700 xu<- max(trapmat[,1])+ 20
701 yu<- max(trapmat[,2])+ 20

```

Finally, use the MCMC code for the binomial encounter model with the clog-log link (`SCR0binom.c1`) and run 5000 iterations. This should take approximately 25 minutes (in real life we would of course run the algorithm a lot longer



705 but for demonstration purposes let's stick with a number of iterations that can  
 706 be run in a manageable amount of time).

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

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

```
714 chain<-mcmc(mod0)
```

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

### 716 Markov chain time series plots

717 Start by looking at time series plots of your Markov chains using `plot(chain)`.  
 718 This command produces a time series plot and marginal posterior density plots  
 719 for each monitored parameter, similar to what we did before using the `hist()`  
 720 and `plot()` commands. Fig. 1.5 shows an example of these plots for  $\sigma$  and  
 721  $\lambda_0$ . Time series plots will tell you several things: First, recall from sec. 1.2.2  
 722 that the way the chains move through the parameter space gives you an idea  
 723 of whether your MH steps are well tuned. If chains were constant over many  
 724 iterations you would need to decrease the tuning parameter of the (Normal)  
 725 proposal distribution. If a chain moves along some gradient to a stationary  
 726 state very slowly, you may want to increase the tuning parameter so that the  
 727 parameter space is explored more efficiently.

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

### 732 1.4.3 Posterior density plots

733 The `plot()` command also produces posterior density plots and it is worthwhile  
 734 to look at those carefully. For parameters with priors that have bounds (e.g.  
 735 Uniform over some interval), you will be able to see if your choice of the prior  
 736 is truncating the posterior distribution. In the context of SCR models, this  
 737 will mostly involve our choice of  $M$ , the size of the augmented data set. If the  
 738 posterior of  $N$  has a lot of mass concentrated close to  $M$  (or equivalently the  
 739 posterior of  $\psi$  has a lot of mass concentrated close to 1), as in the example in  
 740 Fig. 1.6, we have to re-run the analysis with a larger  $M$ . A diffuse posterior  
 741 plot suggests that the parameter may not be well-identified. There may not be  
 742 enough information in your data to estimate model parameters and you may  
 743 have to consider a simpler model. Finally, posterior density plots will show you  
 744 if the posterior distribution is symmetrical or skewed – if the distribution has

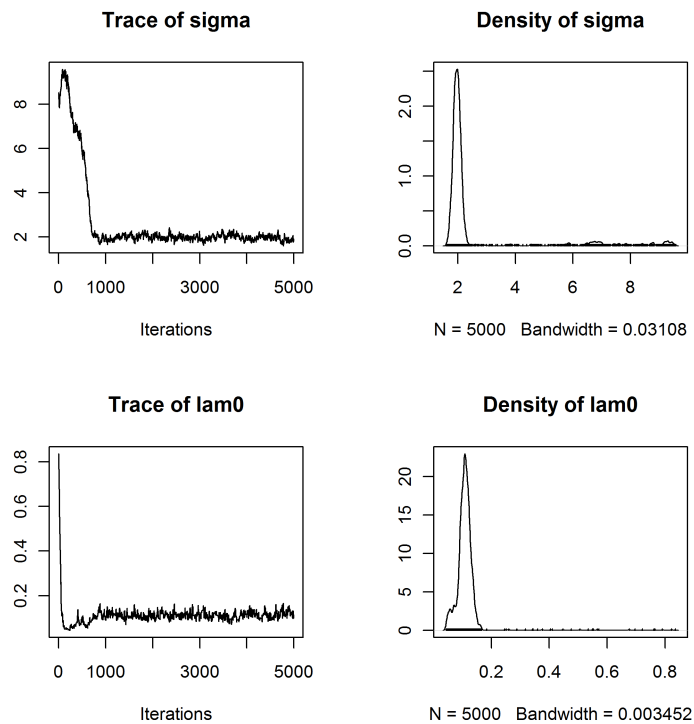


Figure 1.5: Time series and posterior density plots for  $\sigma$  and  $\lambda_0$  for the Fort Drum black bear data.

745 a heavy tail, using the mean as a point estimate of your parameter of interest  
 746 may be biased and you may want to opt for the median or mode instead.

#### 747 1.4.4 Serial autocorrelation and effective sample size

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

```
755 effectiveSize(window(chain, start=1001))
756     sigma     lam0     psi      N
757 93.89807 163.72311 51.96443 46.45394
```

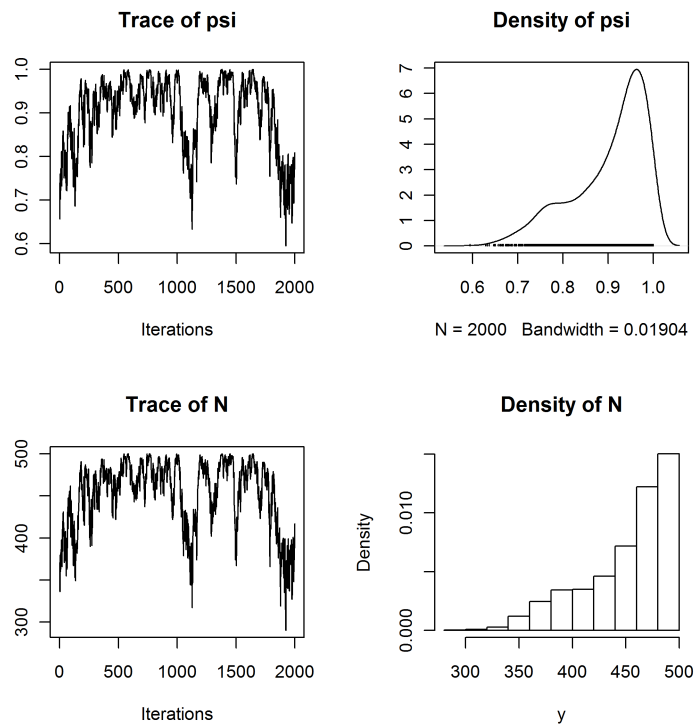


Figure 1.6: Time series and posterior density plots of  $\psi$  and  $N$  for the Fort Drum black bear data truncated by the upper limit of  $M$  (500).

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  $\psi$  and our 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

773 samples to look at the output.

### 774 1.4.5 Summary results

775 Now that we checked that our chains apparently have converged and pretending  
 776 that we have generated enough samples from the posterior distribution, we can  
 777 look at the actual parameter estimates. The `summary()` function will return two  
 778 sets of results: the mean parameter estimates, with their standard deviation,  
 779 the naïve standard error – i.e. your regular standard error calculated for  $T$  (=   
 780 number of iterations) samples without accounting for serial autocorrelation –  
 781 and the Time-series SE (in **WinBUGS** and earlier in this book referred to as  
 782 MC error), which accounts for autocorrelation. Remember our rule of thumb  
 783 that this error decreases with increasing chain length and should be 1% or  
 784 less of the parameter estimate. In **WinBUGS** the MC error is only given in  
 785 the log output within **BUGS** itself. You should adjust the `summary()` call by  
 786 removing the burn-in from calculating parameter summary statistics. To do so,  
 787 use the `window()` command, which lets you specify at which iteration to start  
 788 'counting'. In contrast to **WinBUGS**, which requires you to set the burn-in  
 789 length before you run the model, this command gives us full flexibility to make  
 790 decisions about the burn-in after we have seen the trajectories of our Markov  
 791 chains. For our example, `summary(window(chain, start=1001))` returns the  
 792 following output:

```
793 Iterations = 1001:5000
794 Thinning interval = 1
795 Number of chains = 1
796 Sample size per chain = 4000
797
798 1. Empirical mean and standard deviation for each variable,
799    plus standard error of the mean:
800
801      Mean      SD Naive SE Time-series SE
802 sigma  1.9697  0.12534 0.0019818      0.012792
803 lam0   0.1124  0.01521 0.0002405      0.001311
804 psi    0.7295  0.11794 0.0018648      0.015278
805 N      510.9190 81.99868 1.2965130     10.580567
806
807 2. Quantiles for each variable:
808
809      2.5%      25%      50%      75%      97.5%
810 sigma  1.7288   1.8831   1.9666   2.0517   2.2240
811 lam0   0.0863   0.1008   0.1112   0.1217   0.1449
812 psi    0.5100   0.6423   0.7261   0.8170   0.9549
813 N      359.0000 451.0000 508.0000 572.0000 668.0000
```

814 Looking at the MC errors (column labeled **Time-series SE**), we see that in  
 815 spite of the high autocorrelation, the MC error for  $\sigma$  is below the 1% threshold,  
 816 whereas for all other parameters, MC errors are still above, another indication  
 817 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 km<sup>2</sup>) and we can use summary statistics of the resulting distribution to characterize  $D$ :

```
summary(window(chain[,4]/ 3032.719, start=1001))
```

```
Iterations = 1001:5000
Thinning interval = 1
Number of chains = 1
Sample size per chain = 4000
```

1. Empirical mean and standard deviation for each variable,  
plus standard error of the mean:

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

2. Quantiles for each variable:

	2.5%	25%	50%	75%	97.5%
	0.1184	0.1487	0.1675	0.1886	0.2203

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

#### 1.4.6 Other useful commands

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

```
rejectionRate(chain)
```

	sigma	lam0	psi	N
	0.42988598	0.78775755	0.00000000	0.03160632

Recall (sec. 1.2.2) that rejection rates should lie between 0.2 and 0.8, so our tuning seems to have been appropriate here. Draws of the parameter  $\psi$  are never rejected since we update it with Gibbs sampling, where all candidate values are kept. And since  $N$  is the sum of all  $z_i$ , all it takes for  $N$  to change from one iteration to the next are small changes in the  $z$ -vector, so the rejection rate of  $N$  is always low. If you have run several parallel chains, you can combine them into a single mcmc object using the `mcmc.list()` command on the individual chains (note that each chain has to be converted to an mcmc object before combining them with `mcmc.list()`). You can then easily obtain the Gelman-Rubin diagnostic (Gelman et al., 2004), in **WinBUGS** called Rhat, using `gelman.diag()`,

which will indicate if all chains have converged to the same stationary distribution. For details on these and other functions, see the **coda** manual, which can be found (together with the package) on the CRAN mirror.

## 1.5 Manipulating the state-space

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

As mentioned before, this situation is difficult to handle in **WinBUGS**. In some simple cases we can adjust the state space by setting one of the coordinates 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 space<sup>3</sup>. 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 our **WinBUGS** 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])
```

$\text{In}$  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  $[\text{SX}, \text{SY}]$  pair will be rejected. If  $\text{Pin} = 1$ , this bit of the likelihood is equal to 1, and whether or not 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  $\text{SX}$  and  $\text{SY}$  that it requires.

<sup>3</sup>This idea was pitched to us by Mike Meredith, from WCS Malaysia

In **R**, we are much more flexible, as we can use the actual state-space polygon to constrain  $\mathbf{s}_i$ . 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.

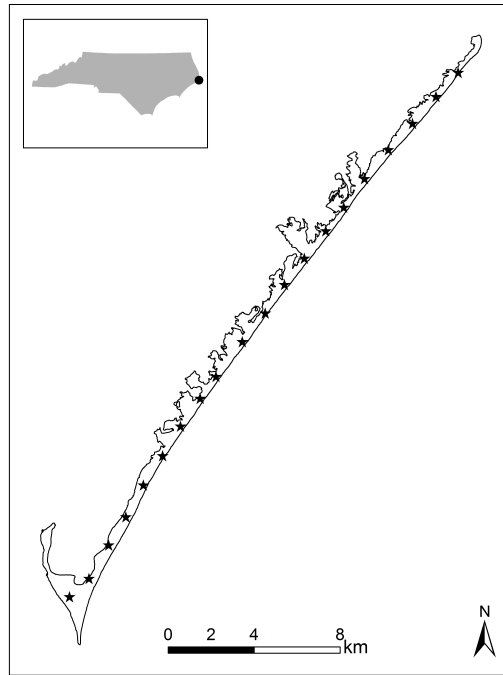


Figure 1.7: Camera traps (stars) set up on South Core Banks, a barrier island within Cape Lookout National Seashore, North Carolina (inset map) to estimate the raccoon population (see Chapt. ?? for details).

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 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 shape files – 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:

```

924 xl= xmin-b
925 xu= xmax+b
926 yl= ymin-b
927 yu= ymax+b
928
929         #create data frame with coordinate pairs
930 dd <- data.frame(Id=c(1,1,1,1,1),X=c(xl,xu,xu,xl,xl),
931   Y=c(yl,yl,yu,yu,yl))
932 ddTable <- data.frame(Id=c(1),Name=c("Item1"))
933         #convert to shapefile, type polygon
934 ddShapefile <- convert.to.shapefile(dd, ddTable, "Id", 5)
935         # name and save to location of choice
936 write.shapefile(ddShapefile, 'c:/Test', arcgis=T)

```

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

```

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

```

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

```

953         #draw candidate value
954 Scand <- as.matrix(cbind(rnorm(M, S[,1], 2), rnorm(M, S[,2], 2)))

```

---

<sup>4</sup>Remember from the previous chapter (??) that the `over` function takes as its second argument (among others) an object of the class “SpatialPolygons” or “SpatialPolygons-DataFrame”. 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.



```

955     #convert to spatial points on UTM (m) scale
956     Scoord<-SpatialPoints(Scand*1000)
957     # check if scand is within the polygon
958     SinPoly<-over(Scoord,SSp)
959
960     for(i in 1:M) {
961         #if scand falls within polygon, continue update
962         if(is.na(SinPoly[i])==FALSE) {
963             ... [rest of the updating step remains the same]

```

964 Note that it is much more time-efficient to draw all  $M$  candidate values for  
965  $s$  and check once if they fall within the state-space, rather than running the  
966 `over()` command for every individual pair of coordinates. To make sure that  
967 our initial values for  $s$  also fall within the polygon of  $\mathcal{S}$ , we use the function  
968 `runifpoint()` from the package `spatstat` (Baddeley and Turner, 2005), which  
969 generates random uniform points within a specified polygon. You'll find this  
970 modified MCMC algorithm (`SCR0poisSSp`) in the **R** package `scrbook`.

971 Finally, observe that we are converting candidate coordinates of  $\mathcal{S}$  back to  
972 meters to match the UTM polygon. In all previous examples, for both the trap  
973 locations and the activity centers we have used UTM coordinates divided by  
974 1000 to estimate  $\sigma$  on a km scale. This is adequate for wide ranging species like  
975 bears. In other cases you may center all coordinates on 0. No matter what kind  
976 of transformation you use on your coordinates, make sure to always convert  
977 candidate values for  $\mathcal{S}$  back to the original scale (UTM) before running the  
978 `over()` command.

## 979 1.6 Increasing computational speed

980 Using custom written MCMC algorithms in **R** is not only more flexible but  
981 can also be faster than using programs such as **JAGS** and especially **BUGS**.  
982 Also, **R** tends to use much less memory than **JAGS**, which can be crucial if  
983 you are running a large model but only have limited memory available. For  
984 example, you will see in Chapt. ?? that even with a reasonable sized data set  
985 certain parameterizations of SCR models can max out the memory of a 16 GB  
986 computer when using **JAGS**. These are mostly the models that require us to  
987 look at individual sampling occasions instead of joining observations for a given  
988 sampling location across the entire study, which requires us to introduce another  
989 for-loop into the **JAGS** model. **BUGS** is limited in the amount of memory it  
990 can access and thus will likely not max out your memory, but as a trade-off,  
991 it will take a long time to run such models. In this chapter we have provided  
992 you with the guidelines to write your own MCMC sampler. But beyond the  
993 material that we have covered there are a number of ways you can make your  
994 sampler more efficient, through parallel computing or by accessing an alterna-  
995 tive computer language such as **c++**. Exploring these options exhaustively is  
996 beyond the scope of this book; instead, in this section we will give you some  
997 pointers to get started with these more advanced computational issues.

### 1.6.1 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 at the same time 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. 1.4.2) 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):

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

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:

```

1042 sfInit( parallel=TRUE, cpus=3 ) #initialize cluster
1043 sfLibrary(scrbook) #export library scrbook
1044 sfExportAll() #export all data in current workspace
1045 sfClusterSetupRNG() #set up random number generator
1046 outL=sfLapply(1:3,wrapper) # execute 'wrapper' 3 times

```

1047 The object outL is a list of length 3, with one out matrix from the function  
 1048 `SCR0binom.c1` for each chain. After computation is complete, terminate the  
 1049 cluster using the command `sfStop()`. Note that the intermediate output of  
 1050 current values and acceptance rates in the **R** console is suppressed when using  
 1051 parallel computing. We can now look at the output as described previously  
 1052 using the package coda, by first defining outL to be a list of mcmc objects.

```

1053 library(coda)
1054 #turn output into MCMC list
1055 res<-mcmc.list(as.mcmc(outL[[1]]),as.mcmc(outL[[2]]),as.mcmc(outL[[3]]))
1056 summary(window(res, start=1001)) #remove first 1000 iterations as burn-in
1057
1058 [... some output removed ...]

```

```

1059
1060      Mean      SD Naive SE Time-series SE
1061 sigma  1.9723  0.13093 0.0011952      0.0087055
1062 lam0   0.1115  0.01535 0.0001401      0.0009003
1063 psi    0.7130  0.10787 0.0009847      0.0077910
1064 N      499.6166 74.74934 0.6823650      5.4232653

```

```

1065
1066 2. Quantiles for each variable:

```

```

1067
1068      2.5%      25%      50%      75%      97.5%
1069 sigma  1.74339  1.8811  1.9637  2.0530  2.2618
1070 lam0   0.08443  0.1007  0.1105  0.1211  0.1438
1071 psi    0.52046  0.6350  0.7093  0.7814  0.9627
1072 N      366.00000 446.0000 497.0000 547.0000 674.0000

```

1073 Now that we have parallel chains we can also use the function `gelman.diag`  
 1074 to evaluate if chains have converged:

```

1075 gelman.diag(window(res, start=1001)) #assess chain convergence

```

```

1076
1077 Potential scale reduction factors:

```

```

1078
1079      Point est. Upper C.I.
1080 sigma      1.01      1.04
1081 lam0       1.01      1.02
1082 psi        1.07      1.21
1083 N          1.07      1.21

```

```

1084
1085 Multivariate psrf

```

```

1086
1087 1.05

```

We can see that estimates are similar to what we observed when running a single chain (see sec. 1.4.2) 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  $\psi$  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:

```
effectiveSize(window(res, start=1001))

      sigma      lam0      psi      N
272.6935 411.8384 167.4192 168.3355
```

## 1.6.2 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 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:

```

1138 gx<-seq(1,10,1)
1139 gy<-seq(1,10,1)
1140 X<-as.matrix(expand.grid(gx, gy))
1141 M<-700
1142 J<-dim(X)[1]
1143 b<-3
1144 xl<-min(gx)-b
1145 xu<-max(gx)+b
1146 yl<-min(gy)-b
1147 yu<-max(gy)+b
1148 S<-cbind(runif(M, xl, xu), runif(M, yl,yu))

```

Next, we can write a “pedestrian” version of **e2dist** and check how long it takes to calculate the squared distance matrix:

```

1151 Dfun<-function(M, J, S, X){
1152   D2<-matrix(0, nrow=M, ncol=J)
1153   for (i in 1:M){
1154     for(j in 1:J){
1155       D2[i,j]<-(S[i,1]-X[j,1])^2 + (S[i,2]-X[j,2])^2
1156     }
1157   }
1158   return(D2)
1159 }
1160 system.time(
1161   (D2R<-Dfun(M, J, S, X))
1162 )
1163
1164 user  system elapsed
1165  0.81    0.01    0.82

```

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

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

---

```

### 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**

```

1180 system.time(
1181   (out<-DfunArma(M,J,S,X)))
1182
1183      user  system elapsed
1184         0        0        0

```

1185 While speed differences of less than 1 second may seem negligible, remember  
 1186 that each command has to be executed at each iteration of the Markov chain. Espe-  
 1187 cially with time-consuming models such as those for open populations (Chapt.  
 1188 ??) or multi-session models (Chapt. ??) we believe that **C++** holds large  
 1189 potential to make implementation of such models more feasible.

## 1.7 Summary and Outlook

In a nutshell, programs like **WinBUGS** do everything that we went through in this chapter (and quite a bit more). Looking through your model, they determine which parameters they can use standard Gibbs sampling for (i.e. for conjugate full conditional distributions). Then, they determine whether to use adaptive rejection sampling, slice sampling or – in the ‘worst’ case – Metropolis-Hastings sampling for the other full conditionals (how the sampler is chosen differs among softwares). For MH sampling, they will automatically tune the updater so that it works efficiently.

Although these programs are flexible and extremely useful to perform MCMC simulations, it sometimes is more efficient to develop your own MCMC algorithm. Building an MCMC code follows three basic steps: Identify your model including priors and express full conditional distributions for each model parameter. If full conditionals are parametric distributions, use Gibbs sampling to draw candidate parameter values from those distributions; otherwise use Metropolis-Hastings sampling to draw candidate values from a proposal distribution and accept or reject them based on their posterior probability densities.

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

Another situations in which using **BUGS/JAGS** becomes increasingly complicated or inefficient is when using point processes other than the homogeneous Binomial point process (“uniformity of density”) which underlies the basic SCR model (see sec. ?? in Chapt. ??). In Chapt. ?? you already saw an example of an inhomogeneous point process model and we briefly introduce a different point processes, implemented using a custom-made MCMC algorithm, in Chapt. ?. Finally, as mentioned earlier, the following chapters, ?? and ??, deal with unmarked or partially marked populations using hand-made MCMC algorithms to handle the (partially) latent individual encounter histories. While some of these models can be written in **BUGS/JAGS**, they are painstakingly slow; others (for example the classes of models considered in Chapt. ??) cannot be implemented in **BUGS/JAGS** at all and we have to either use likelihood based inference or develop our own MCMC algorithms. In conclusion, while you can certainly get by using **BUGS/JAGS** for standard SCR models, knowing how to write your own MCMC sampler allows you 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.

- 1266 Gopalaswamy, A. M., Royle, A. J., Hines, J., Singh, P., Jathanna, D., Ku-  
1267 mar, N. S., and Karanth, K. U. (2012), “Program SPACECAP: software for  
1268 estimating animal density using spatially explicit capturerecapture models,”  
1269 *Methods in Ecology and Evolution*, online early, r package version 1.0.4.
- 1270 Hastings, W. (1970), “Monte Carlo sampling methods using Markov chains and  
1271 their applications,” *Biometrika*, 57, 97–109.
- 1272 Knaus, J. (2010), *snowfall: Easier cluster computing (based on snow)*., r package  
1273 version 1.84.
- 1274 Lewin-Koh, N. J., Bivand, R., contributions by Edzer J. Pebesma, Archer, E.,  
1275 Baddeley, A., Bibiko, H.-J., Dray, S., Forrest, D., Friendly, M., Giraudoux, P.,  
1276 Golicher, D., Rubio, V. G., Hausmann, P., Hufthammer, K. O., Jagger, T.,  
1277 Luque, S. P., MacQueen, D., Niccolai, A., Short, T., Stabler, B., and Turner,  
1278 R. (2011), *maptools: Tools for reading and handling spatial objects*, r package  
1279 version 0.8-10.
- 1280 Link, W. A. and Barker, R. J. (2010), *Bayesian Inference: With Ecological*  
1281 *Applications*, London, UK: Academic Press.
- 1282 Metropolis, N., Rosenbluth, A., Rosenbluth, M., Teller, A., Teller, E., et al.  
1283 (1953), “Equation of state calculations by fast computing machines,” *The*  
1284 *journal of chemical physics*, 21, 1087–1092.
- 1285 Metropolis, N. and Ulam, S. (1949), “The Monte Carlo method,” *Journal of the*  
1286 *American Statistical Association*, 44, 335–341.
- 1287 Neal, R. (2003), “Slice sampling,” *Annals of Statistics*, 31, 705–741.
- 1288 Pebesma, E. and Bivand, R. (2011), *Package ‘sp’*, r package version 0.9-91.
- 1289 Plummer, M. (2003), “JAGS: A program for analysis of Bayesian graphical mod-  
1290 els using Gibbs sampling,” in *Proceedings of the 3rd International Workshop*  
1291 *on Distributed Statistical Computing (DSC 2003)*. March, pp. 20–22.
- 1292 Plummer, M., Best, N., Cowles, K., and Vines, K. (2006), “CODA: Convergence  
1293 Diagnosis and Output Analysis for MCMC,” *R News*, 6, 7–11.
- 1294 Robert, C. P. and Casella, G. (2004), *Monte Carlo statistical methods*, New  
1295 York, USA: Springer.
- 1296 — (2010), *Introducing Monte Carlo Methods with R*, New York, USA: Springer.
- 1297 Roberts, G. O. and Rosenthal, J. S. (1998), “Optimal scaling of discrete ap-  
1298 proximations to Langevin diffusions,” *Journal of the Royal Statistical Society:*  
1299 *Series B (Statistical Methodology)*, 60, 255–268.
- 1300 Sklyar, O., Murdoch, D., Smith, M., Eddelbuettel, D., and François, R. (2010),  
1301 *inline: Inline C, C++, Fortran function calls from R*, r package version 0.3.8.

- 1302 Sollmann, R., Gardner, B., Parsons, A., Stocking, J., McClintock, B., Simons,  
1303 T., Pollock, K., and O'Connell, A. (2013), "A spatial mark-resight model  
1304 augmented with telemetry data," *Ecology*.
- 1305 Stabler, B. (2006), *shapefiles: Read and Write ESRI Shapefiles*, r package version  
1306 0.6.
- 1307 Tierney, L., Rossini, A. J., Li, N., and Sevcikova, H. (2011), *snow: Simple*  
1308 *Network of Workstations*, r package version 0.3-7.