

HMSC-R 3.0: Getting started with HMSC-R: univariate models

Gleb Tikhonov, Oystein H. Opedal, Nerea Abrego, Aleksi Lehikoinen & Otso Ovaskainen

11 March 2019

Introduction

The Hierarchical Modelling of Species Communities (HMSC) framework is a statistical framework for analysis of multivariate data, typically from species communities. Here, we demonstrate how to get started with HMSC. While HMSC is primarily meant for multivariate data, this vignette illustrates the basics with univariate models.

To get HMSC in use, you need to first install it [<https://github.com/hmsc-r/HMSC>] and then load it.

```
library(Hmsc)
```

We set the random number seed to make the results presented here reproducible.

```
set.seed(1)
```

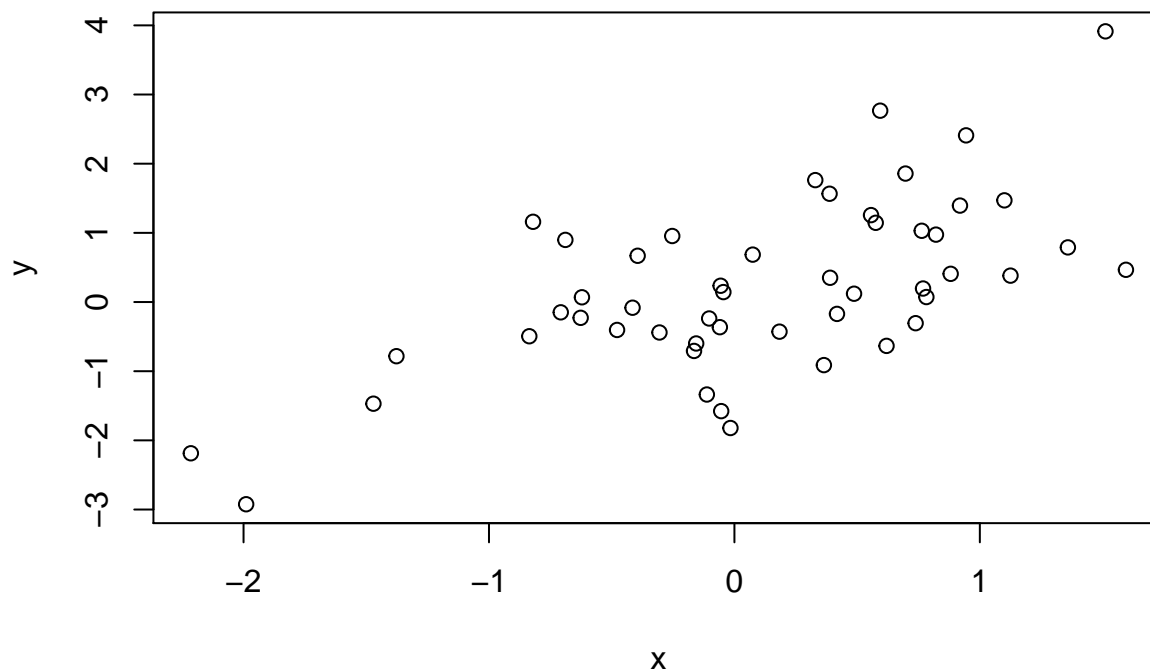
Linear model

As the first case study, we use HMSC to fit a univariate linear model. To relate the results to something that we expect the reader to be familiar with, we also apply the basic `lm` function to the same data and compare the results.

Generating simulated data

For illustrative purposes, we use simulated data for which we know the parameter values.

```
n = 50
x = rnorm(n)
alpha = 0
beta = 1
sigma = 1
L = alpha + beta*x
y = L + rnorm(n, sd = sigma)
plot(x,y)
```



Here n is the number of datapoints, x is a continuous covariate, α and β are the true parameters for intercept and slope, L is the linear predictor, and y is the response variable. We note that the data are simulated by the standard linear model with normally distributed residuals, $y_i = \alpha + \beta x_i + \epsilon_i$, where $\epsilon_i \sim N(0, \sigma^2)$.

Fitting models and looking at parameter estimates

The standard way to analyze these kinds of data with the maximum likelihood inference is to use the `lm` function:

```
da = data.frame(x,y)
m.lm = lm(y ~ x, data=da)
summary(m.lm)
```

```
##
## Call:
## lm(formula = y ~ x, data = da)
##
## Residuals:
```

	Min	1Q	Median	3Q	Max
	-1.92760	-0.66898	-0.00225	0.48768	2.34858

```
##
## Coefficients:
```

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	0.1219	0.1394	0.875	0.386
x	0.9545	0.1681	5.679	7.73e-07 ***

```
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.9781 on 48 degrees of freedom
## Multiple R-squared:  0.4019, Adjusted R-squared:  0.3894
## F-statistic: 32.25 on 1 and 48 DF,  p-value: 7.726e-07
```

We note that, as expected, the parameter estimates roughly correspond to the values we assumed for the intercept ($\alpha = 0$) and slope ($\beta = 1$) when generating the data.

To conduct the analogous analyses with HMSC, we first construct the model as

```
Y=as.matrix(y)
XData = data.frame(x=x)
m = Hmsc(Y=Y, XData=XData, XFormula=~x)
```

While the `lm` function constructs the model and fits it at the same time, the `Hmsc` function only constructs the model object. To fit the HMSC model with Bayesian inference, we use the `sampleMcmc` function.

```
thin = 2
nChains = 2
samples = 500
transient = 250
verbose = 250
m = sampleMcmc(m, thin = thin, samples = samples, transient = transient,
               nChains = nChains, verbose = verbose)
```

```
## [1] "Computing chain 1"
## [1] "Chain 1, iteration 250 of 1250, (transient)"
## [1] "Chain 1, iteration 500 of 1250, (sampling)"
## [1] "Chain 1, iteration 750 of 1250, (sampling)"
## [1] "Chain 1, iteration 1000 of 1250, (sampling)"
## [1] "Chain 1, iteration 1250 of 1250, (sampling)"
## [1] "Computing chain 2"
## [1] "Chain 2, iteration 250 of 1250, (transient)"
## [1] "Chain 2, iteration 500 of 1250, (sampling)"
## [1] "Chain 2, iteration 750 of 1250, (sampling)"
## [1] "Chain 2, iteration 1000 of 1250, (sampling)"
## [1] "Chain 2, iteration 1250 of 1250, (sampling)"
```

Here we have selected some parameters (`thin`, `samples`, `transient`, `nChains`) that control how the parameters are estimated (more precisely, how the posterior distribution is sampled). We will ignore these for a moment but return to their meaning soon. The main thing is that now the model object `m` also includes estimated parameters, in the same way as the object `m.lm` includes the parameters estimated by the `lm` function.

To see the parameter estimates, we may e.g. extract the posterior distribution from the model object and convert it into a coda object.

```
mpost = convertToCodaObject(m)
summary(mpost$Beta)
```

```
##
## Iterations = 252:1250
## Thinning interval = 2
## Number of chains = 2
## Sample size per chain = 500
##
## 1. Empirical mean and standard deviation for each variable,
```

```
## plus standard error of the mean:
##
##               Mean      SD Naive SE Time-series SE
## B[(Intercept) (C1), sp1 (S1)] 0.1239 0.1447 0.004577      0.004578
## B[x (C2), sp1 (S1)]          0.9466 0.1691 0.005348      0.005519
##
## 2. Quantiles for each variable:
##
##               2.5%      25%      50%      75%      97.5%
## B[(Intercept) (C1), sp1 (S1)] -0.1601 0.02675 0.1223 0.226 0.3879
## B[x (C2), sp1 (S1)]          0.6277 0.82681 0.9482 1.057 1.2689
```

To assess model fit in terms of R^2 , we apply the `evaluateModelFit` function to posterior predictive distribution computed by the function `computePredictedValues`.

```
preds = computePredictedValues(m)
evaluateModelFit(hM=m, predY=preds)
```

```
## $RMSE
## [1] 0.9583804
##
## $R2
## [1] 0.4018712
```

We note that the parameter estimates and R^2 given by HMSC are highly consisted with those given by the `lm` function. We however note that the two approaches are not identical as `lm` applies the maximum likelihood framework (and thus e.g. yields confidence intervals) whereas HMSC applies the Bayesian framework (and thus e.g. assumes prior distributions and yields credible intervals).

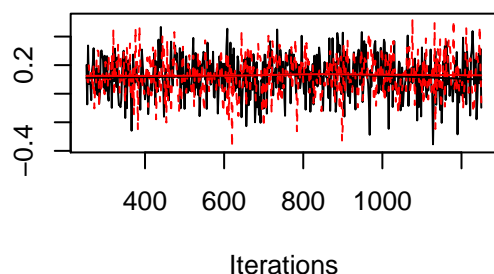
Model fit can be evaluated in many ways. The other measure of model fit returned here by the `evaluateModelFit` function is RMSE, i.e. the root-mean-square error.

Checking MCMC diagnostics

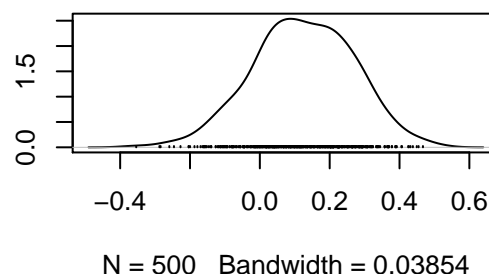
Let us then return to the parameters that guide posterior sampling in the MCMC algorithm. We first plot the trace-plots of the β -parameters.

```
plot(mpost$Beta)
```

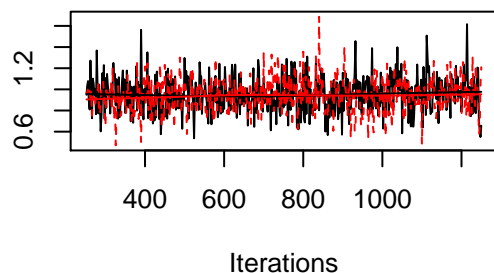
Trace of B[(Intercept) (C1), sp1 (S1)]



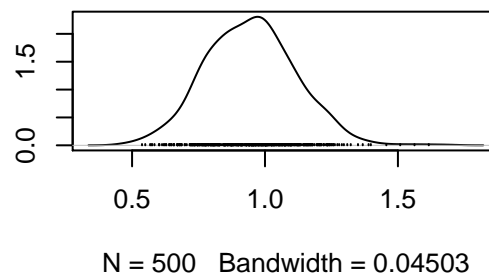
Density of B[(Intercept) (C1), sp1 (S1)]



Trace of B[x (C2), sp1 (S1)]



Density of B[x (C2), sp1 (S1)]



The black and red colors show the two (`nChains=2`) independent MCMC chains. Both chains start from iteration 251, as by `transient=250` we have chosen to ignore iterations before that as a possible transient. The chains have run in total 1250 iterations each, as we selected to obtain 500 samples, and we have recorded only every 2:th step (`thin=2`) of the iterations.

These trace-plots look essentially as good as they can ever look like. First of all, the two chains yield essentially identical results. Second, the chains mix very well, i.e. they go fast up and down without any apparent autocorrelation. Third, they seemed to have reached a stationary distribution, as e.g. the first half of the recorded iterations looks statistically identical to the second half of the recorded iterations.

We may evaluate MCMC convergence also more quantitatively in terms of effective sample size and potential scale reduction factor.

```
effectiveSize(mpost$Beta)
```

```
## B[(Intercept) (C1), sp1 (S1)]      B[x (C2), sp1 (S1)]
##                               1000.0000                945.0366
```

```
gelman.diag(mpost$Beta,multivariate=FALSE)$psrf
```

```
##                               Point est. Upper C.I.
## B[(Intercept) (C1), sp1 (S1)]  1.0042620  1.0215579
## B[x (C2), sp1 (S1)]           0.9990046  0.9995102
```

We observe that the effective sample sizes are very close to the theoretical value of the actual number of samples, which is 1000 (500 per chain). This indicates that there is very little autocorrelation among consecutive samples. The potential scale reduction factors are very close to one, which indicates that two chains gave consistent results, as suggested by visual inspection of the trace plots.

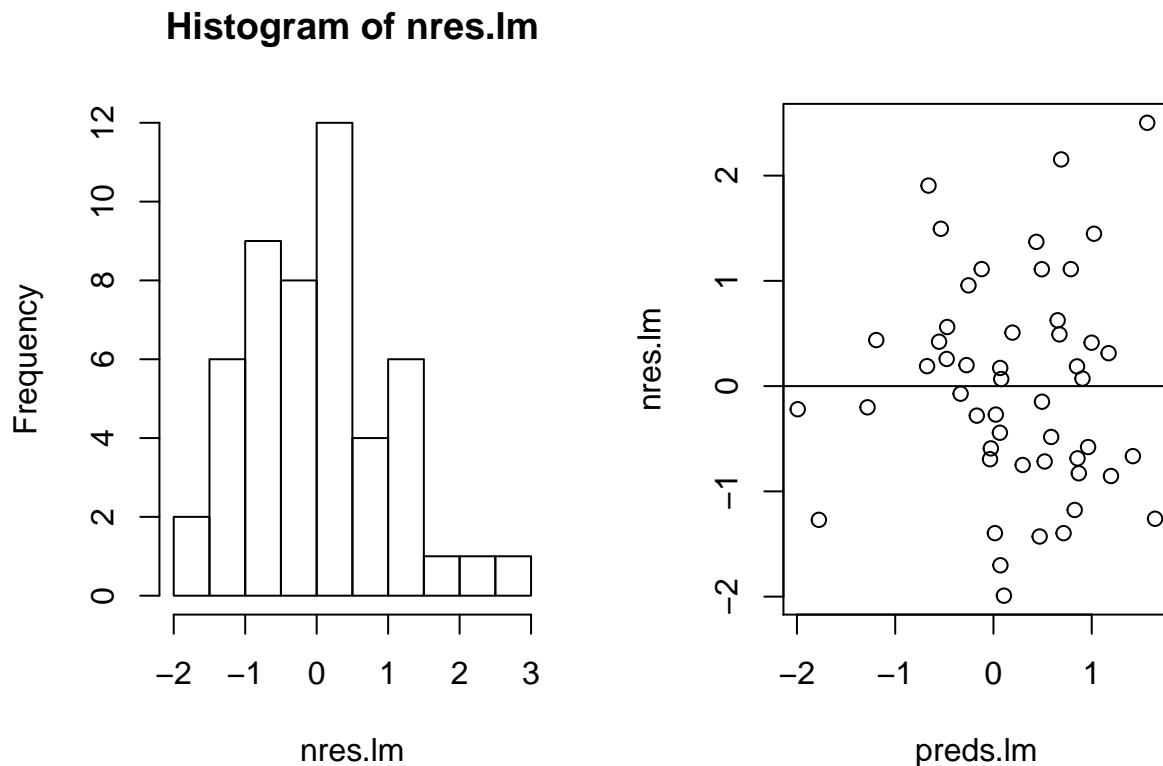
In summary, the MCMC diagnostics did not indicate any problems with MCMC convergence. This means that

the posterior sample is likely to be representative of the true posterior distribution, and thus the inference from the model can be trusted. If the MCMC convergence would have indicated problems, we should have refitted the model using different parameters. If the trace-plots would have suggested a present of a transient (the early iterations would have looked different from the later iterations), we should have increased the amount of transient (also called burn-in) iterations to be discarded. If the chains would have shown autocorrelations and/or differed from each other, we should have increased the number of iterations. This can be done by increasing either the number of samples, or by keeping the number of samples fixed but increasing thinning. We recommend the latter, as increasing the number of samples can make the model objects very big and lead to computationally expensive postprocessing. If one needs to run a million iterations, we recommend doing so by `samples=1000` and `thin=1000` rather than `samples=1000000` and `thin=1`. In our view, 1000 samples are typically sufficient to evaluate e.g. posterior means and credible intervals with sufficient accuracy.

Checking the assumptions of the linear model

We finally note that when applying the linear model, it is a good practice to examine if the assumptions of the linear model are valid. In the context of the `lm` function, this can be done e.g. by constructing the following diagnostic plots.

```
nres.lm = rstandard(m.lm)
preds.lm = fitted.values(m.lm)
par(mfrow=c(1,2))
hist(nres.lm)
plot(preds.lm,nres.lm)
abline(a=0,b=0)
```

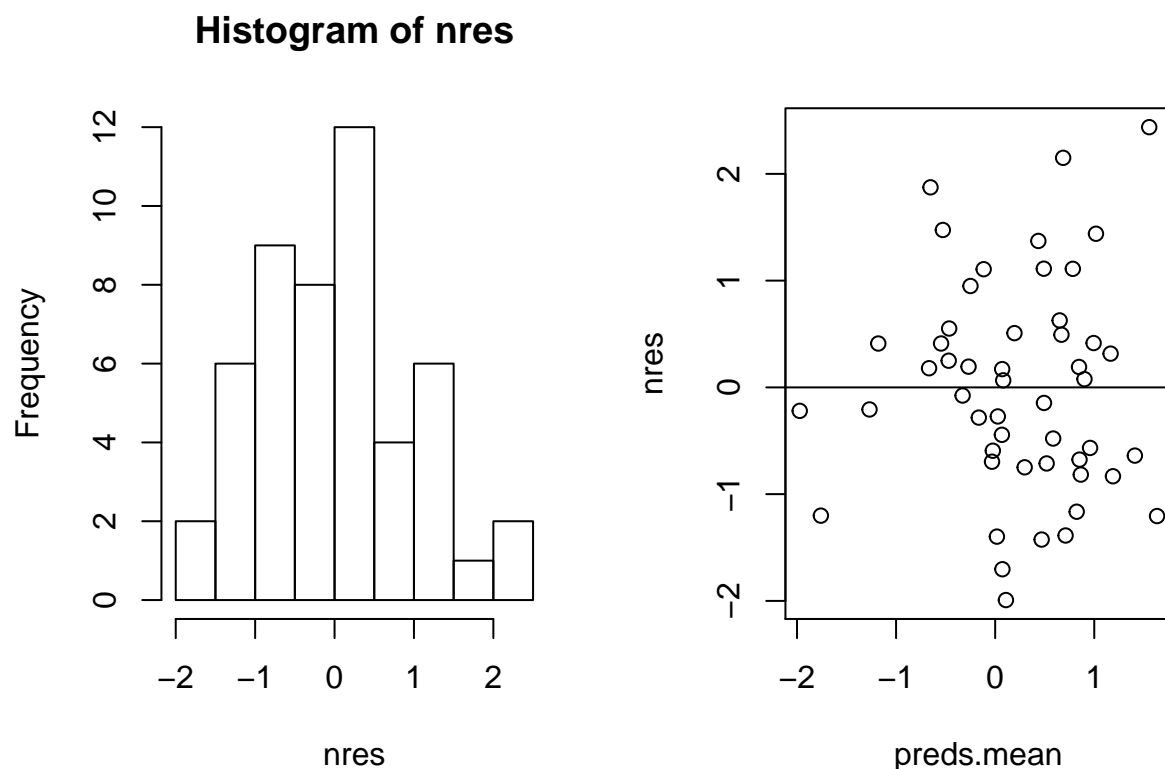


The first plot shows that the residuals conform well to the assumption of normality, and the second plot shows

that the residuals are homoscedastic. This is not surprising, as the data were simulated from the linear model.

Structural model assumptions can be checked with HMSC basically in the same way as with the standard linear model. However, while model fitted with the `lm` function simply typing `plot(m.lm)` would provide many diagnostic plots “automatically”, with HMSC there is no such built-in functionality. This is because typical applications of HMSC relate to much more complex models where it is not straightforward to decide what a “standard” diagnostic plot should be. To generate diagnostic plots with HMSC, we first summarize the posterior predictive distribution into posterior mean and then extract and standardize the residuals.

```
preds.mean = apply(preds, FUN=mean, MARGIN=1)
nres = scale(y-preds.mean)
par(mfrow=c(1,2))
hist(nres)
plot(preds.mean,nres)
abline(a=0,b=0)
```



These diagnostic plots are essentially identical to those obtained for the linear model. This is to be expected, as we observed earlier that the parameter estimates (and hence the fitted values and residuals) were consistent between the two approaches.

Generalized linear models

Above, we have fitted a linear model with normally distributed residuals. As this has been set as the default option for HMSC, we defined our model above simply as

```
m = Hmsc(Y=Y, XData=XData, XFormula=~x)
```

whereas the full version that makes the assumption of normality transparent would be

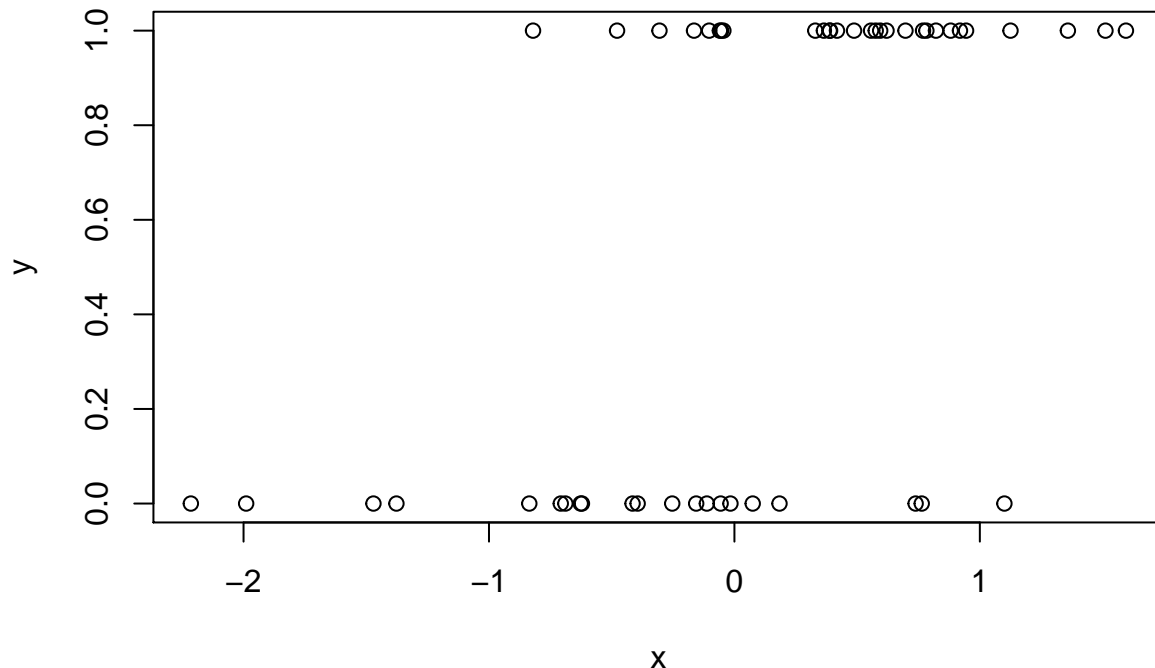
```
m = Hmsc(Y=Y, XData=XData, XFormula=~x, distr = "normal")
```

In addition to the normal distribution, HMSC-R 3.0 allows three other types of link functions and error distributions: probit model for presence-absence data, and Poisson and log-normal Poisson models for count data. We note that while also many other kinds of link functions and error distributions can be quite straightforwardly implemented in the univariate framework, this is not the case for HMSC due to its multivariate and hierarchical nature. Thus, implementing other kinds of link functions and error distributions is a challenge for the future.

Probit model for presence-absence data

A common case with community ecology data is that of presence-absence data, meaning that the response variable is either 0 (the species is absent) or 1 (the species is present). In HMSC-R, such data can be modelled with probit regression. For those readers not familiar with probit regression, we note that it is similar to logistic regression, it just applies a slightly different link-function (probit instead of logit). We next repeat the above exercise but do so for data that conform the assumptions of probit regression.

```
y = 1*(L+ rnorm(n, sd = 1)>0)
plot(x,y)
```



```
Y=as.matrix(y)
m = Hmsc(Y=Y, XData=XData, XFormula=~x, distr="probit")
```



```
m = sampleMcmc(m, thin = thin, samples = samples, transient = transient,
               nChains = nChains, verbose = verbose)
```

```
## [1] "Computing chain 1"
## [1] "Chain 1, iteration 250 of 1250, (transient)"
## [1] "Chain 1, iteration 500 of 1250, (sampling)"
## [1] "Chain 1, iteration 750 of 1250, (sampling)"
## [1] "Chain 1, iteration 1000 of 1250, (sampling)"
## [1] "Chain 1, iteration 1250 of 1250, (sampling)"
## [1] "Computing chain 2"
## [1] "Chain 2, iteration 250 of 1250, (transient)"
## [1] "Chain 2, iteration 500 of 1250, (sampling)"
## [1] "Chain 2, iteration 750 of 1250, (sampling)"
## [1] "Chain 2, iteration 1000 of 1250, (sampling)"
## [1] "Chain 2, iteration 1250 of 1250, (sampling)"
```

```
mpost = convertToCodaObject(m)
summary(mpost$Beta)
```

```
##
## Iterations = 252:1250
## Thinning interval = 2
## Number of chains = 2
## Sample size per chain = 500
##
## 1. Empirical mean and standard deviation for each variable,
##    plus standard error of the mean:
##
##              Mean      SD Naive SE Time-series SE
## B[(Intercept) (C1), sp1 (S1)] 0.1078 0.2023 0.006396      0.008367
## B[x (C2), sp1 (S1)]          1.1292 0.3275 0.010357      0.017248
##
## 2. Quantiles for each variable:
##
##              2.5%      25%      50%      75%      97.5%
## B[(Intercept) (C1), sp1 (S1)] -0.3115 -0.02189 0.1158 0.2452 0.4781
## B[x (C2), sp1 (S1)]          0.5695  0.88769 1.1174 1.3354 1.7985
```

```
effectiveSize(mpost$Beta)
```

```
## B[(Intercept) (C1), sp1 (S1)]          B[x (C2), sp1 (S1)]
##              615.0609                  448.3885
```

```
gelman.diag(mpost$Beta,multivariate=FALSE)$psrf
```

```
##              Point est. Upper C.I.
## B[(Intercept) (C1), sp1 (S1)]  1.000012  1.003602
## B[x (C2), sp1 (S1)]          1.004896  1.007883
```

```
preds = computePredictedValues(m)
evaluateModelFit(hM=m, predY=preds)
```

```
## $RMSE
## [1] 0.4064478
##
## $AUC
## [1] 0.816092
```

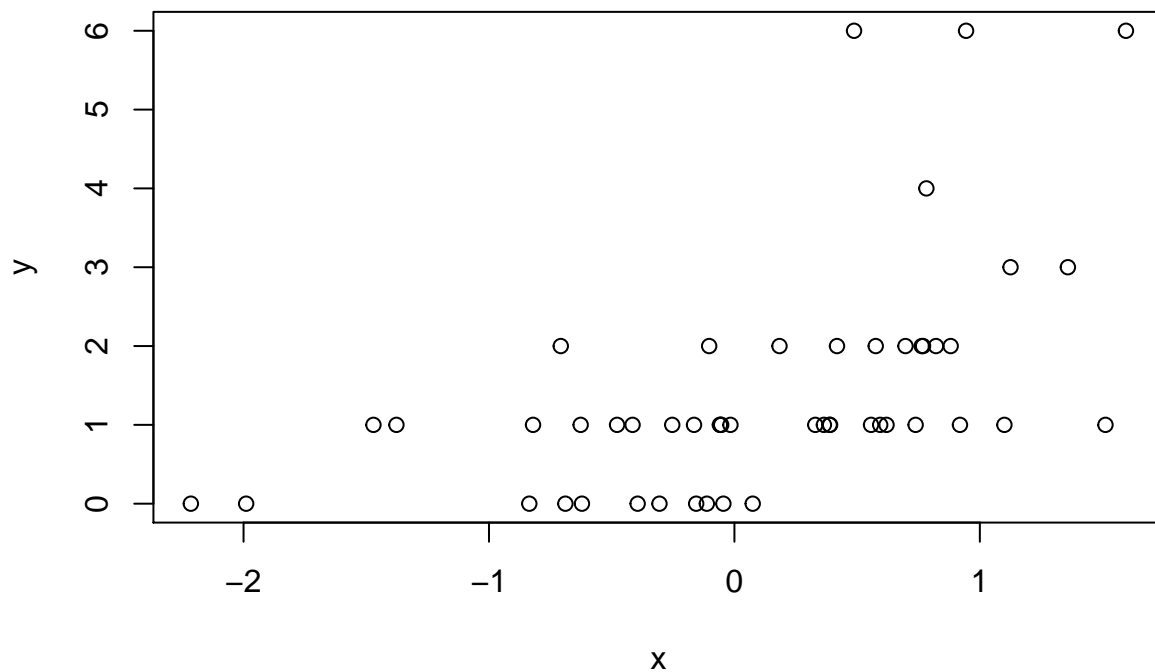
```
##
## $TjurR2
## [1] 0.3010591
```

Compared to the case of the linear model, we observe that the effective sample size is somewhat smaller (achieving MCMC convergence is generally more challenging in non-normal models), and the parameter estimates include more uncertainty (0/1 data are less informative than normally distributed data). We further note that instead of the R^2 of the linear model, model fit is now evaluated in terms of AUC and Tjur R^2 .

Poisson model for count data

In another common case with community ecology data the response variable is a non-negative integer 0,1,2,3,..., representing the count of individuals. In HMSC-R, such data can be modelled with Poisson regression. We next repeat the above exercise but do so for data that conform the assumptions of Poisson regression.

```
y = y = rpois(n, lambda = exp(L))
plot(x,y)
```



```
Y=as.matrix(y)
m = Hmsc(Y=Y, XData=XData, XFormula=~x, distr="poisson")
m = sampleMcmc(m, thin = thin, samples = samples, transient = transient,
               nChains = nChains, verbose = verbose)
```

```
## [1] "Computing chain 1"
## [1] "Chain 1, iteration 250 of 1250, (transient)"
## [1] "Chain 1, iteration 500 of 1250, (sampling)"
```

```

## [1] "Chain 1, iteration 750 of 1250, (sampling)"
## [1] "Chain 1, iteration 1000 of 1250, (sampling)"
## [1] "Chain 1, iteration 1250 of 1250, (sampling)"
## [1] "Computing chain 2"
## [1] "Chain 2, iteration 250 of 1250, (transient)"
## [1] "Chain 2, iteration 500 of 1250, (sampling)"
## [1] "Chain 2, iteration 750 of 1250, (sampling)"
## [1] "Chain 2, iteration 1000 of 1250, (sampling)"
## [1] "Chain 2, iteration 1250 of 1250, (sampling)"

mpost = convertToCodaObject(m)
summary(mpost$Beta)

##
## Iterations = 252:1250
## Thinning interval = 2
## Number of chains = 2
## Sample size per chain = 500
##
## 1. Empirical mean and standard deviation for each variable,
##    plus standard error of the mean:
##
##              Mean      SD Naive SE Time-series SE
## B[(Intercept) (C1), sp1 (S1)] 0.02475 0.09792 0.003096      0.02411
## B[x (C2), sp1 (S1)]          0.87227 0.15252 0.004823      0.04964
##
## 2. Quantiles for each variable:
##
##              2.5%      25%      50%      75%  97.5%
## B[(Intercept) (C1), sp1 (S1)] -0.1989 -0.03208 0.03408 0.09467 0.1865
## B[x (C2), sp1 (S1)]          0.6156  0.76442 0.85449 0.96911 1.1918

effectiveSize(mpost$Beta)

## B[(Intercept) (C1), sp1 (S1)]          B[x (C2), sp1 (S1)]
##              20.89455                  11.42770

gelman.diag(mpost$Beta,multivariate=FALSE)$psrf

##              Point est. Upper C.I.
## B[(Intercept) (C1), sp1 (S1)]    1.329767    2.198235
## B[x (C2), sp1 (S1)]              1.120191    1.219361

preds = computePredictedValues(m, expected = FALSE)
evaluateModelFit(hM=m, predY=preds)

## $RMSE
## [1] 1.30384
##
## $SR2
## [1] 0.2941907
##
## $O.AUC
## [1] 0.8368298
##
## $O.TjurR2
## [1] 0.2535641

```

```
##
## $O.RMSE
## [1] 0.3770874
##
## $C.SR2
## [1] 0.07663744
##
## $C.RMSE
## [1] 1.450879
```

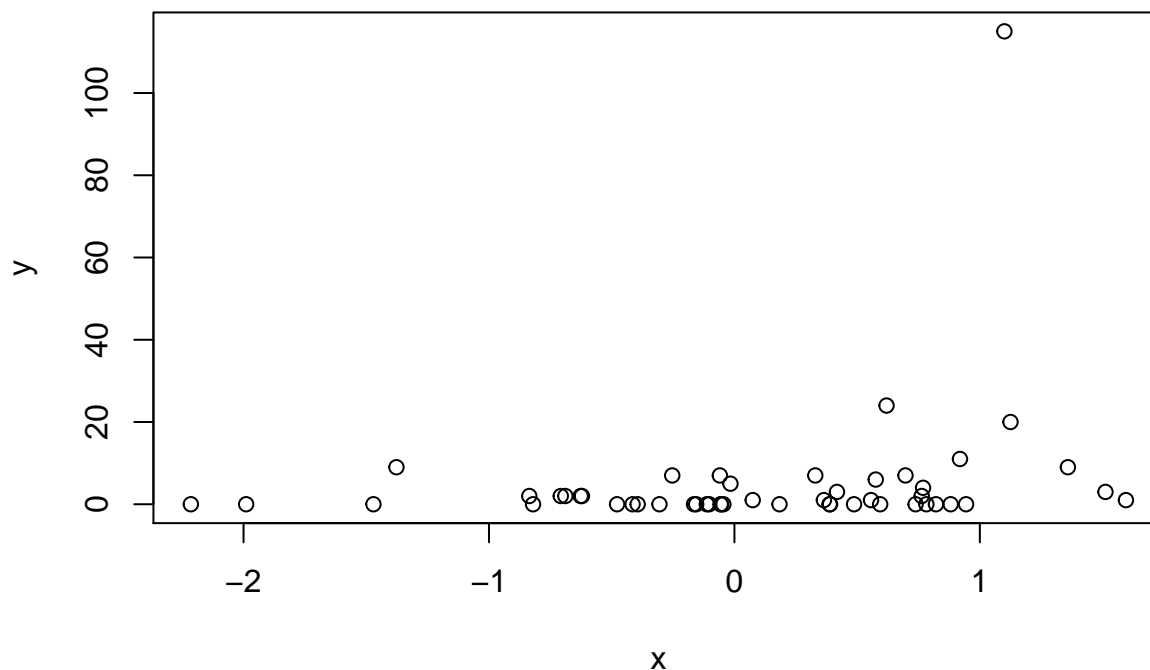
In this case, the MCMC diagnostics are much worse than for the probit model. Unfortunately, achieving MCMC convergence for Poisson model in the HMSC context is highly challenging. To get results that can be trusted, one should increase the thinning and thus run the model longer.

For the Poisson model, many kinds of measures of model fit are given. The measure **SR2** can be viewed as a pseudo- R^2 . Whereas the R^2 of a linear model can be computed as the squared pearson correlation between predicted and true values, we have computed **SR2** as the squared spearman correlation between observed and predicted values. As the Poisson model predicts counts, it can be used to separate whether a species is present (count>0) or absent (count=0). For evaluating how well species occurrences (indicated by **O.**) are predicted, the same measures (**AUC** and Tjur R^2) can be computed as for the probit model. Sometimes it is also be of interest to examine how well the model is able to predict abundance variation in cases where the species is present, and thus ignoring absences. This is done with the measures indicated by **C.**, where **C** refers to “conditional on presence”. Note that here we have generated the predictions with the option **expected = FALSE**. In this case, the predictions are not expected values (e.g. on average we expect to see 2.3 individuals), but a posterior predictive distribution of data (i.e., actual counts 0,1,2,3,... that involve the Poisson error). We have selected **expected = FALSE** so that presence-absences can be inferred from the predictions, enabling us to compute also the **O.** and **C.** variants of measures of model fit.

Log-normal Poisson model for count data

The Poisson model is often not a good model for ecological count data as it does not allow sufficient amount of variation around the expectation. To allow for more variation, HMSC-R includes the possibility to fit a log-normal Poisson model. We note that another standard option would be to fit the Negative Binomial model, but this model is not included currently in HMSC-R. We next repeat the above exercise but do so for data that conform the assumptions of log-normal Poisson regression.

```
y = rpois(n, lambda = exp(L+rnorm(n, sd=2)))
plot(x,y)
```



```
Y=as.matrix(y)
m = Hmsc(Y=Y, XData=XData, XFormula=~x, distr="lognormal poisson")
m = sampleMcmc(m, thin = thin, samples = samples, transient = transient,
               nChains = nChains, verbose = verbose)
```

```
## [1] "Computing chain 1"
## [1] "Chain 1, iteration 250 of 1250, (transient)"
## [1] "Chain 1, iteration 500 of 1250, (sampling)"
## [1] "Chain 1, iteration 750 of 1250, (sampling)"
## [1] "Chain 1, iteration 1000 of 1250, (sampling)"
## [1] "Chain 1, iteration 1250 of 1250, (sampling)"
## [1] "Computing chain 2"
## [1] "Chain 2, iteration 250 of 1250, (transient)"
## [1] "Chain 2, iteration 500 of 1250, (sampling)"
## [1] "Chain 2, iteration 750 of 1250, (sampling)"
## [1] "Chain 2, iteration 1000 of 1250, (sampling)"
## [1] "Chain 2, iteration 1250 of 1250, (sampling)"
```

```
mpost = convertToCodaObject(m)
summary(mpost$Beta)
```

```
##
## Iterations = 252:1250
## Thinning interval = 2
## Number of chains = 2
## Sample size per chain = 500
##
```

```
## 1. Empirical mean and standard deviation for each variable,
##    plus standard error of the mean:
##
##              Mean      SD Naive SE Time-series SE
## B[(Intercept) (C1), sp1 (S1)] -0.2394 0.3671 0.01161      0.06924
## B[x (C2), sp1 (S1)]          0.7750 0.3537 0.01119      0.03490
##
## 2. Quantiles for each variable:
##
##              2.5%      25%      50%      75% 97.5%
## B[(Intercept) (C1), sp1 (S1)] -1.10127 -0.4508 -0.1955 0.02432 0.3816
## B[x (C2), sp1 (S1)]          0.07397 0.5540 0.7645 0.99615 1.5101

effectiveSize(mpost$Beta)

## B[(Intercept) (C1), sp1 (S1)]          B[x (C2), sp1 (S1)]
##              126.7203                      374.8118

gelman.diag(mpost$Beta,multivariate=FALSE)$psrf

##              Point est. Upper C.I.
## B[(Intercept) (C1), sp1 (S1)] 1.313909 2.106846
## B[x (C2), sp1 (S1)]          1.049157 1.161965

preds = computePredictedValues(m, expected = FALSE)
evaluateModelFit(hM=m, predY=preds)

## $RMSE
## [1] 16.78035
##
## $SR2
## [1] 0.06281139
##
## $O.AUC
## [1] 0.6344
##
## $O.TjurR2
## [1] 0.05536
##
## $O.RMSE
## [1] 0.4884734
##
## $C.SR2
## [1] 0.03956341
##
## $C.RMSE
## [1] 23.44627
```

We note that the output is similar to that given for the standard Poisson model.

Hurdle models

Ecological count data are often dominated by zeros, i.e. they are zero-inflated. One standard model that can be fitted to such data is the Zero-Inflated Poisson model. While HMSC-R 3.0 does not include this or other zero-inflated models, we note that it is always possible to apply the closely-related Hurdle approach, i.e. to analyze separately occurrence data $1*(Y>0)$ by probit regression, and abundance conditional on presence

$Y[Y==0] = NA$ by a suitable model of abundance. We note that these two aspects of the data (occurrence and abundance conditional on presence) are statistically independent of each other, and thus it is interesting to compare results obtained for them, e.g. to see if variation in occurrence and variation in abundance are explained by the same environmental covariates.

How to select the best model?

When fitting generalized linear models with the maximum likelihood framework, a common way of comparing models is to use AIC. For example, in case of count data, AIC could be used to compare the Poisson model, the Negative Binomial model, and the Zero-Inflated Poisson model. In case of HMSC-R, model selection is recommended to be conducted through evaluating predictive performance through cross-validation. We illustrate cross-validation strategies below in the context of mixed models.

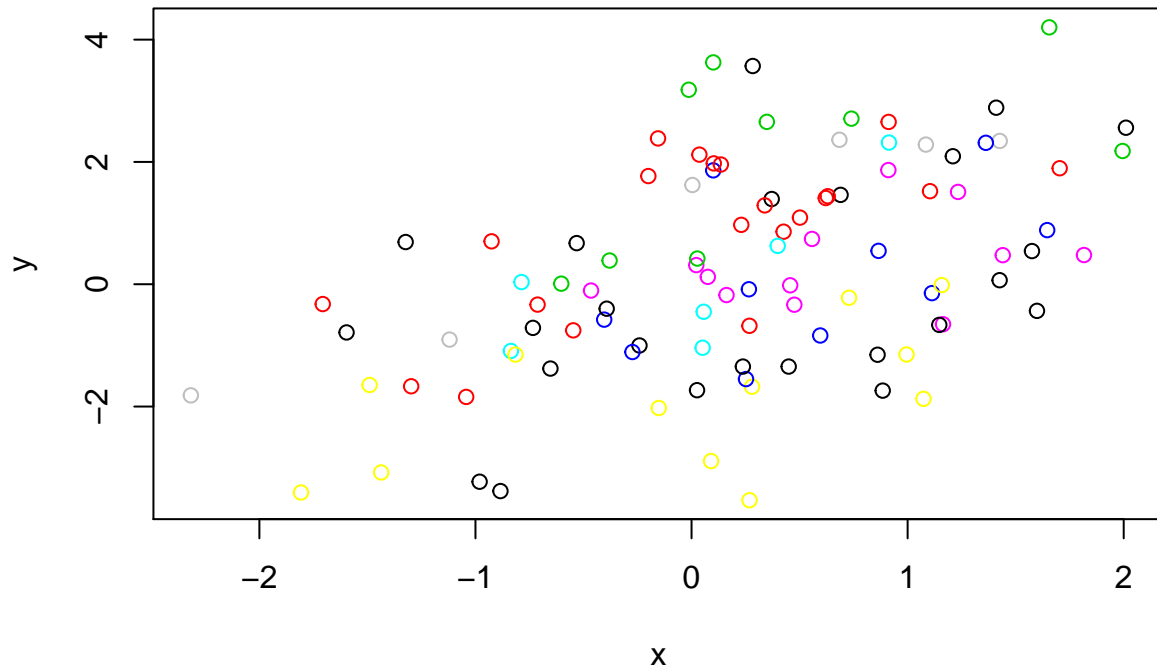
Mixed models: fixed and random effects

We next turn to mixed models, i.e. models with both fixed and random effects. For simplicity, we consider here linear models only, even if similar analyses could be done with probit, Poisson or log-normal Poisson models.

Hierarchical structure: sampling units within plots

We first consider a hierarchical study design in which 100 sampling units belong to a discrete set of `np=10` plots. We assume that the plots have an additive effect to the response variable, i.e. we consider a random intercept model.

```
n = 100
x = rnorm(n)
alpha = 0
beta = 1
sigma = 1
L = alpha + beta*x
np = 10
sigma.plot = 1
sample.id = 1:n
plot.id = sample(1:np,n,replace = TRUE)
ap = rnorm(np, sd = sigma.plot)
a = ap[plot.id]
y = L + a + rnorm(n, sd = sigma)
plot.id = as.factor(plot.id)
plot(x,y,col=plot.id)
```



```
XData = data.frame(x=x)
Y=as.matrix(y)
```

We note that in the maximum likelihood framework, one option for fitting a mixed model would be to apply the `lmer` function as `lmer(y~x+(1|plot.id),data=da)`. Fitting a similar model with HMSC can be done as follows.

```
studyDesign = data.frame(sample = as.factor(sample.id), plot = as.factor(plot.id))
rL = HmscRandomLevel(units = studyDesign$plot)
rL$nfMin = 1
rL$nfMax = 1
m = Hmsc(Y=Y, XData=XData, XFormula=~x,
         studyDesign=studyDesign, ranLevels=list("plot"=rL))
m = sampleMcmc(m, thin = thin, samples = samples, transient = transient,
              nChains = nChains, verbose = verbose)
```

```
## [1] "Computing chain 1"
## [1] "Chain 1, iteration 250 of 1250, (transient)"
## [1] "Chain 1, iteration 500 of 1250, (sampling)"
## [1] "Chain 1, iteration 750 of 1250, (sampling)"
## [1] "Chain 1, iteration 1000 of 1250, (sampling)"
## [1] "Chain 1, iteration 1250 of 1250, (sampling)"
## [1] "Computing chain 2"
## [1] "Chain 2, iteration 250 of 1250, (transient)"
## [1] "Chain 2, iteration 500 of 1250, (sampling)"
## [1] "Chain 2, iteration 750 of 1250, (sampling)"
## [1] "Chain 2, iteration 1000 of 1250, (sampling)"
```



```
## [1] "Chain 2, iteration 1250 of 1250, (sampling)"
```

Here we have included two new input parameters. First, `studyDesign` defines the nature of the study design, here the individual sampling units and the plots to which the sampling units belong to. Second, `ranLevels` includes a list of those aspects of the study design to which a random effect will be included, here the plot. The object `rL` was created with the function `HmscRandomLevel` assuming the default options except that with `rL$nfMin = 1` and `rL$nfMax = 1` we have fixed the number of the underlying latent factors to one. This is because we assume a univariate model, and thus it is not meaningful to have several latent variables. We note that in usual applications of HMSC, the data are multivariate, and the latent factor approach is used to reduce the dimensionality by having a smaller number of latent factors than the number of species.

Let us then evaluate model fit as we have done above.

```
preds = computePredictedValues(m)
MF = evaluateModelFit(hM=m, predY=preds)
MF$R2
```

```
## [1] 0.7251942
```

What is predicted here are the same data that were used to fit the model, and thus R^2 measures explanatory power rather than predictive power. If we would add more predictors, we could make R^2 approach to one even if the predictors would represent just random noise, which phenomenon is known as overfitting. For this reason, the predictive power of the model should be evaluated more critically by cross-validation.

Examples of cross-validation strategies

To apply cross-validation, we first assign the samples randomly into a number of folds. For example, we may apply two-fold cross-validation across the samples by making the following partition

```
partition = createPartition(m, nfolds = 2, column = "sample")
partition
```

```
## [1] 1 2 1 2 2 1 1 1 1 1 2 2 2 2 2 1 1 1 2 1 1 1 1 2 2 2 1 2 1 2 1 2 1
## [36] 2 1 1 2 1 2 2 1 2 2 1 1 1 1 1 2 2 1 2 1 1 1 2 2 1 2 2 1 2 2 1 1 1 1
## [71] 2 1 2 1 2 2 2 2 2 1 2 2 2 2 1 2 2 1 2 1 1 2 2 2 2 2 1 1 2 1
```

The idea of cross-validation is that when making predictions for a particular fold, data from that fold is not used for parameter estimation. Model fitting and predictions are made separately for each fold, so that in the end a prediction is obtained for each sampling unit.

```
preds = computePredictedValues(m, partition=partition)
```

```
## [1] "Cross-validation, fold 1 out of 2"
## [1] "Computing chain 1"
## [1] "Chain 1, iteration 250 of 1250, (transient)"
## [1] "Chain 1, iteration 500 of 1250, (sampling)"
## [1] "Chain 1, iteration 750 of 1250, (sampling)"
## [1] "Chain 1, iteration 1000 of 1250, (sampling)"
## [1] "Chain 1, iteration 1250 of 1250, (sampling)"
## [1] "Computing chain 2"
## [1] "Chain 2, iteration 250 of 1250, (transient)"
## [1] "Chain 2, iteration 500 of 1250, (sampling)"
## [1] "Chain 2, iteration 750 of 1250, (sampling)"
## [1] "Chain 2, iteration 1000 of 1250, (sampling)"
## [1] "Chain 2, iteration 1250 of 1250, (sampling)"
## [1] "Cross-validation, fold 2 out of 2"
## [1] "Computing chain 1"
## [1] "Chain 1, iteration 250 of 1250, (transient)"
```

```
## [1] "Chain 1, iteration 500 of 1250, (sampling)"
## [1] "Chain 1, iteration 750 of 1250, (sampling)"
## [1] "Chain 1, iteration 1000 of 1250, (sampling)"
## [1] "Chain 1, iteration 1250 of 1250, (sampling)"
## [1] "Computing chain 2"
## [1] "Chain 2, iteration 250 of 1250, (transient)"
## [1] "Chain 2, iteration 500 of 1250, (sampling)"
## [1] "Chain 2, iteration 750 of 1250, (sampling)"
## [1] "Chain 2, iteration 1000 of 1250, (sampling)"
## [1] "Chain 2, iteration 1250 of 1250, (sampling)"
```

```
MF = evaluateModelFit(hM=m, predY=preds)
MF$R2
```

```
## [1] 0.6140673
```

As expected, the cross-validation based predictive R^2 is lower than the explanatory R^2 .

Making the two-fold cross-validation required fitting the model twice, and more generally performing k -fold cross-validation requires fitting the model k times, which can become computationally intensive. For this reason, we have applied here two-fold cross-validation rather than e.g. leave-one-out cross-validation. Increasing the number of folds means that more data are available for fitting the model, which can be expected to lead to greater predictive performance. In this sense, the predictive power estimated by two-fold cross-validation is likely to underestimate the true predictive power of the full model fitted to all data.

Cross-validation can be made in many ways, and how exactly it should be made depends on which aspect of predictive power one wishes to measure. To illustrate, let us partition the plots rather than the sampling units into different folds.

```
partition = createPartition(m, nfolds = 2, column = "plot")
t(cbind(plot.id,partition)[1:15,])
```

```
##           [,1] [,2] [,3] [,4] [,5] [,6] [,7] [,8] [,9] [,10] [,11] [,12]
## plot.id      4  10   4   1   3   2   1  10   5   8   8   8
## partition    2   1   2   1   1   1   1   1   2   2   2   2
##           [,13] [,14] [,15]
## plot.id      9   7   6
## partition    2   1   2
```

As seen by examining the correspondence between plots and the partition (shown here for the first 15 sampling units), now all sampling units belonging to a particular plot have been included in the same fold.

```
preds = computePredictedValues(m,partition=partition)
```

```
## [1] "Cross-validation, fold 1 out of 2"
## [1] "Computing chain 1"
## [1] "Chain 1, iteration 250 of 1250, (transient)"
## [1] "Chain 1, iteration 500 of 1250, (sampling)"
## [1] "Chain 1, iteration 750 of 1250, (sampling)"
## [1] "Chain 1, iteration 1000 of 1250, (sampling)"
## [1] "Chain 1, iteration 1250 of 1250, (sampling)"
## [1] "Computing chain 2"
## [1] "Chain 2, iteration 250 of 1250, (transient)"
## [1] "Chain 2, iteration 500 of 1250, (sampling)"
## [1] "Chain 2, iteration 750 of 1250, (sampling)"
## [1] "Chain 2, iteration 1000 of 1250, (sampling)"
## [1] "Chain 2, iteration 1250 of 1250, (sampling)"
## [1] "Cross-validation, fold 2 out of 2"
```

```
## [1] "Computing chain 1"
## [1] "Chain 1, iteration 250 of 1250, (transient)"
## [1] "Chain 1, iteration 500 of 1250, (sampling)"
## [1] "Chain 1, iteration 750 of 1250, (sampling)"
## [1] "Chain 1, iteration 1000 of 1250, (sampling)"
## [1] "Chain 1, iteration 1250 of 1250, (sampling)"
## [1] "Computing chain 2"
## [1] "Chain 2, iteration 250 of 1250, (transient)"
## [1] "Chain 2, iteration 500 of 1250, (sampling)"
## [1] "Chain 2, iteration 750 of 1250, (sampling)"
## [1] "Chain 2, iteration 1000 of 1250, (sampling)"
## [1] "Chain 2, iteration 1250 of 1250, (sampling)"

MF = evaluateModelFit(hM=m, predY=preds)
MF$R2

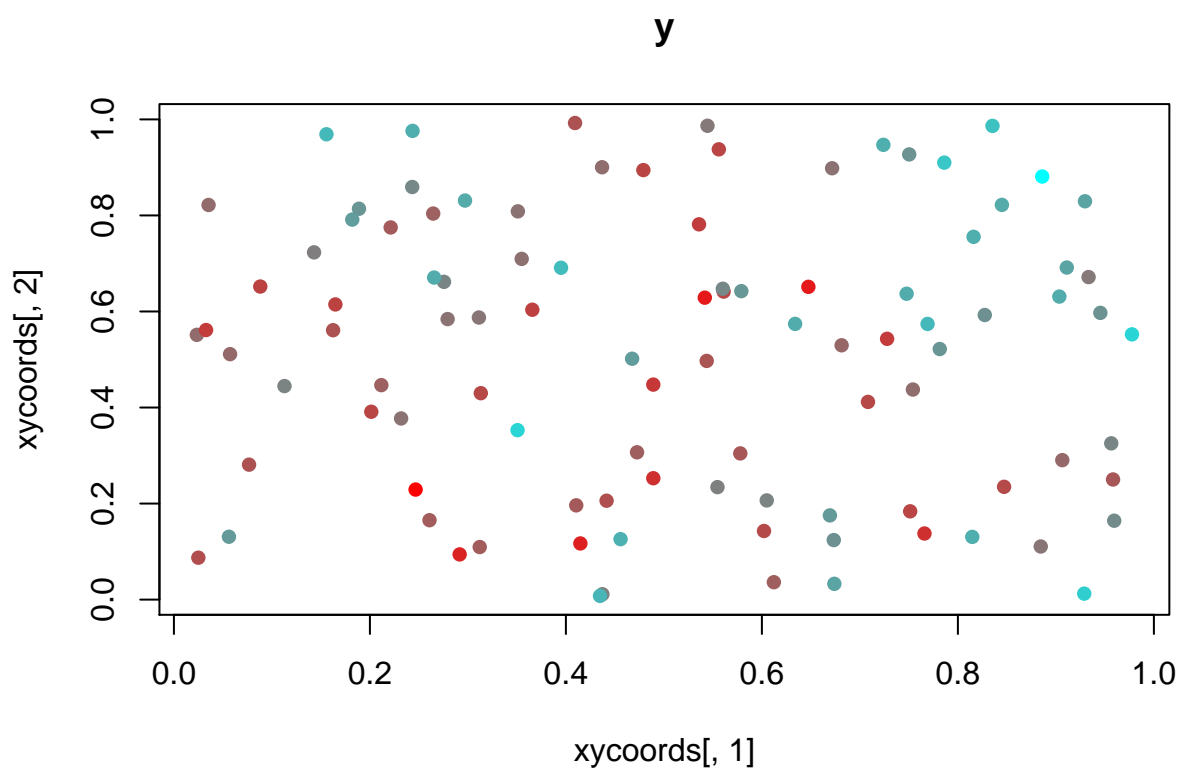
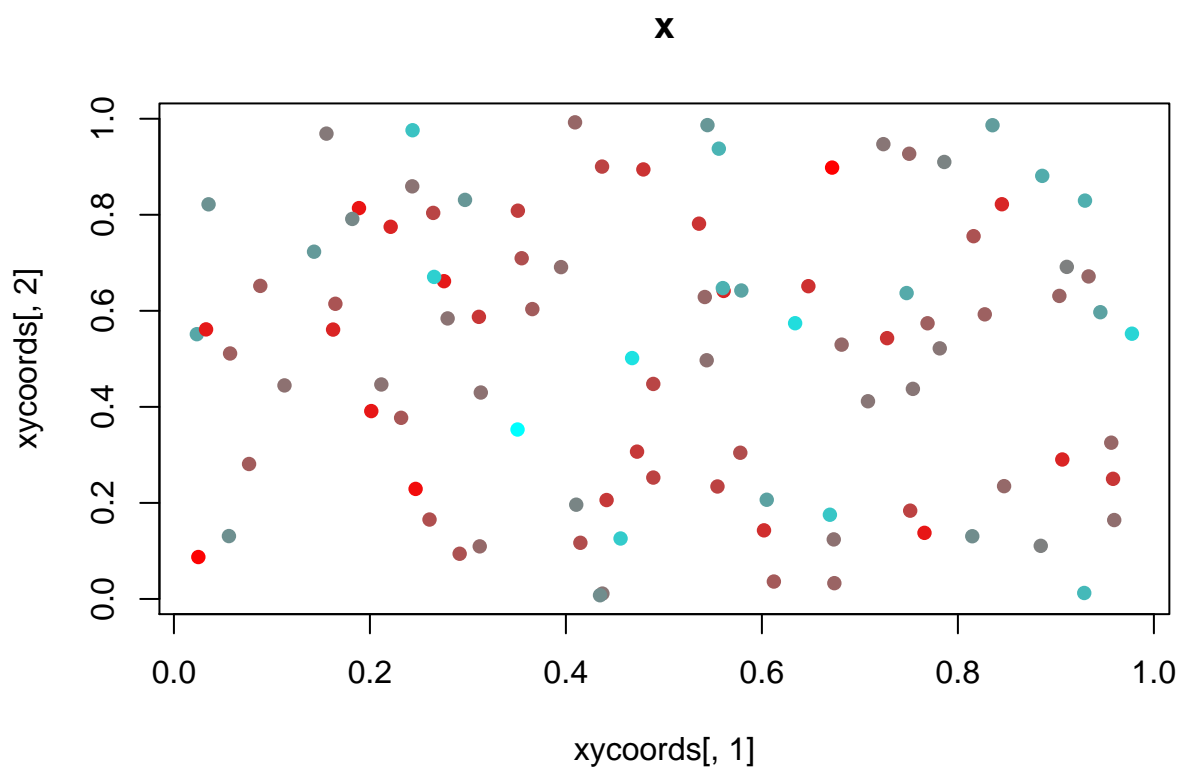
## [1] 0.2178831
```

The predictive power of the model is smaller when assigning entire plots rather than individual sampling units into folds. This is because the prediction task is now more difficult: when making a prediction for a particular sampling unit, the model was now fitted without training data from any other sampling units in the same plot. Thus the model has no possibility to estimate the actual random effect for the focal plot, and thus its predictive power is based solely on the fixed effects.

Spatial structure: including coordinates of sampling units

We next consider spatial data by assuming that each sampling unit is associated with two-dimensional spatial coordinates. We generate the data by assuming a spatially structured residual with exponentially decreasing spatial covariance function.

```
sigma.spatial = 2
alpha.spatial = 0.5
sample.id = rep(NA,n)
for (i in 1:n){
  sample.id[i] = paste0("location_",as.character(i))
}
sample.id = as.factor(sample.id)
xycoords = matrix(runif(2*n), ncol=2)
rownames(xycoords) = sample.id
colnames(xycoords) = c("x-coordinate", "y-coordinate")
a = mvrnorm(mu=rep(0,n),
            Sigma = sigma.spatial*exp(-as.matrix(dist(xycoords))/alpha.spatial))
y = L + a + rnorm(n, sd = sigma)
Y=as.matrix(y)
colfunc = colorRampPalette(c("cyan", "red"))
ncols = 100
cols = colfunc(100)
for (i in 1:2){
  if (i==1) value = x
  if (i==2) value = y
  value = value-min(value)
  value = 1+(ncols-1)*value/max(value)
  plot(xycoords[,1],xycoords[,2],col=cols[value],pch=16,main=c("x","y")[i])
}
```



The spatial structure is visible in the data as nearby sampling units have similar responses y , even if the predictor x is not (due to our assumptions) spatially autocorrelated.

To fit a spatial model with HMSC, we construct the random effect using the `sData` input argument.

```
studyDesign = data.frame(sample = sample.id)
rL = HmscRandomLevel(sData = xycoords)
rL$nfMin = 1
rL$nfMax = 1
m = Hmsc(Y=Y, XData=XData, XFormula=~x,
        studyDesign=studyDesign, ranLevels=list("sample"=rL))
```

Again, we fixed the number of latent factors to one because we consider here a univariate model. Model fitting and evaluation of explanatory and predictive powers can be done as usually.

```
m = sampleMcmc(m, thin = thin, samples = samples, transient = transient,
              nChains = nChains, verbose = verbose)
```

```
## [1] "Computing chain 1"
## [1] "Chain 1, iteration 250 of 1250, (transient)"
## [1] "Chain 1, iteration 500 of 1250, (sampling)"
## [1] "Chain 1, iteration 750 of 1250, (sampling)"
## [1] "Chain 1, iteration 1000 of 1250, (sampling)"
## [1] "Chain 1, iteration 1250 of 1250, (sampling)"
## [1] "Computing chain 2"
## [1] "Chain 2, iteration 250 of 1250, (transient)"
## [1] "Chain 2, iteration 500 of 1250, (sampling)"
## [1] "Chain 2, iteration 750 of 1250, (sampling)"
## [1] "Chain 2, iteration 1000 of 1250, (sampling)"
## [1] "Chain 2, iteration 1250 of 1250, (sampling)"
```

```
preds = computePredictedValues(m)
MF = evaluateModelFit(hM=m, predY=preds)
MF$R2
```

```
## [1] 0.7191066
```

```
partition = createPartition(m, nfolds = 2, column = "sample")
preds = computePredictedValues(m, partition=partition)
```

```
## [1] "Cross-validation, fold 1 out of 2"
## [1] "Computing chain 1"
## [1] "Chain 1, iteration 250 of 1250, (transient)"
## [1] "Chain 1, iteration 500 of 1250, (sampling)"
## [1] "Chain 1, iteration 750 of 1250, (sampling)"
## [1] "Chain 1, iteration 1000 of 1250, (sampling)"
## [1] "Chain 1, iteration 1250 of 1250, (sampling)"
## [1] "Computing chain 2"
## [1] "Chain 2, iteration 250 of 1250, (transient)"
## [1] "Chain 2, iteration 500 of 1250, (sampling)"
## [1] "Chain 2, iteration 750 of 1250, (sampling)"
## [1] "Chain 2, iteration 1000 of 1250, (sampling)"
## [1] "Chain 2, iteration 1250 of 1250, (sampling)"
## [1] "Cross-validation, fold 2 out of 2"
## [1] "Computing chain 1"
## [1] "Chain 1, iteration 250 of 1250, (transient)"
## [1] "Chain 1, iteration 500 of 1250, (sampling)"
## [1] "Chain 1, iteration 750 of 1250, (sampling)"
```

```
## [1] "Chain 1, iteration 1000 of 1250, (sampling)"
## [1] "Chain 1, iteration 1250 of 1250, (sampling)"
## [1] "Computing chain 2"
## [1] "Chain 2, iteration 250 of 1250, (transient)"
## [1] "Chain 2, iteration 500 of 1250, (sampling)"
## [1] "Chain 2, iteration 750 of 1250, (sampling)"
## [1] "Chain 2, iteration 1000 of 1250, (sampling)"
## [1] "Chain 2, iteration 1250 of 1250, (sampling)"

MF = evaluateModelFit(hM=m, predY=preds)
MF$R2
```

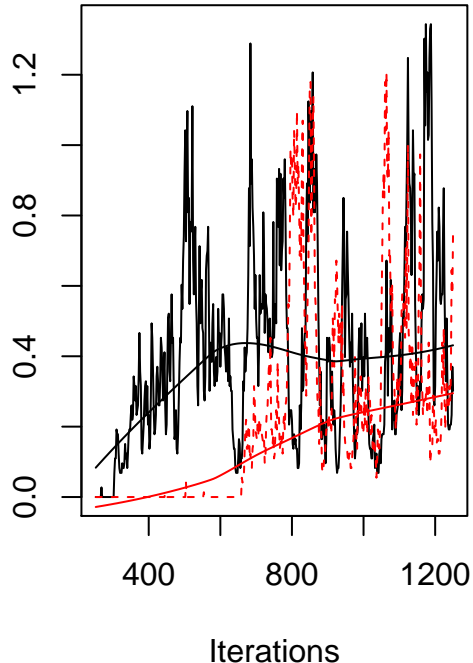
```
## [1] 0.4173968
```

As the model includes a spatially structured random effect, its predictive power is based on both the fixed and the random effects. Concerning the latter, the model can utilize observed data from nearby sampling units included in model fitting when predicting the response for a focal sampling unit that is not included in model fitting.

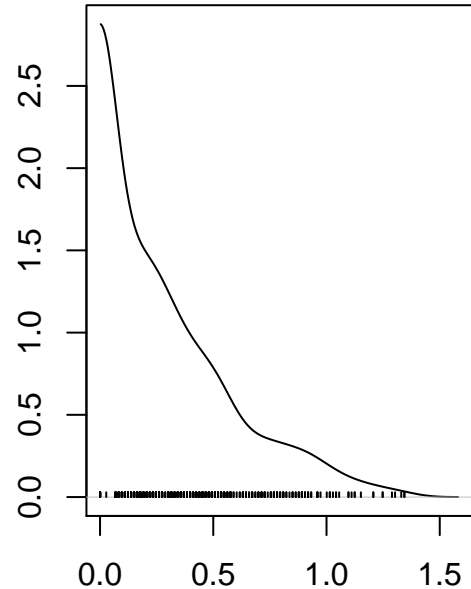
The estimated spatial scale of the estimated random effect is found from the parameter `Alpha[[1]]`, for which we show also the MCMC trace plot.

```
mpost = convertToCodaObject(m)
plot(mpost$Alpha[[1]])
```

Trace of Alpha1[factor1]



Density of Alpha1[factor1]



N = 500 Bandwidth = 0.07901

```
summary(mpost$Alpha[[1]])
```

```
##
## Iterations = 252:1250
```

```
## Thinning interval = 2
## Number of chains = 2
## Sample size per chain = 500
##
## 1. Empirical mean and standard deviation for each variable,
##    plus standard error of the mean:
##
##           Mean           SD      Naive SE Time-series SE
##      0.317204      0.298293      0.009433      0.048552
##
## 2. Quantiles for each variable:
##
##      2.5%      25%      50%      75%      97.5%
## 0.00000 0.08227 0.24682 0.47993 1.04213
```

For comparison, let us fit a non-spatial model to the same data.

```
m = Hmsc(Y=Y, XData=XData, XFormula=~x)
m = sampleMcmc(m, thin = thin, samples = samples, transient = transient,
               nChains = nChains, verbose = verbose)
```

```
## [1] "Computing chain 1"
## [1] "Chain 1, iteration 250 of 1250, (transient)"
## [1] "Chain 1, iteration 500 of 1250, (sampling)"
## [1] "Chain 1, iteration 750 of 1250, (sampling)"
## [1] "Chain 1, iteration 1000 of 1250, (sampling)"
## [1] "Chain 1, iteration 1250 of 1250, (sampling)"
## [1] "Computing chain 2"
## [1] "Chain 2, iteration 250 of 1250, (transient)"
## [1] "Chain 2, iteration 500 of 1250, (sampling)"
## [1] "Chain 2, iteration 750 of 1250, (sampling)"
## [1] "Chain 2, iteration 1000 of 1250, (sampling)"
## [1] "Chain 2, iteration 1250 of 1250, (sampling)"
```

```
preds = computePredictedValues(m)
MF = evaluateModelFit(hM=m, predY=preds)
MF$R2
```

```
## [1] 0.3818968
```

```
preds = computePredictedValues(m,partition=partition)
```

```
## [1] "Cross-validation, fold 1 out of 2"
## [1] "Computing chain 1"
## [1] "Chain 1, iteration 250 of 1250, (transient)"
## [1] "Chain 1, iteration 500 of 1250, (sampling)"
## [1] "Chain 1, iteration 750 of 1250, (sampling)"
## [1] "Chain 1, iteration 1000 of 1250, (sampling)"
## [1] "Chain 1, iteration 1250 of 1250, (sampling)"
## [1] "Computing chain 2"
## [1] "Chain 2, iteration 250 of 1250, (transient)"
## [1] "Chain 2, iteration 500 of 1250, (sampling)"
## [1] "Chain 2, iteration 750 of 1250, (sampling)"
## [1] "Chain 2, iteration 1000 of 1250, (sampling)"
## [1] "Chain 2, iteration 1250 of 1250, (sampling)"
## [1] "Cross-validation, fold 2 out of 2"
## [1] "Computing chain 1"
```

```
## [1] "Chain 1, iteration 250 of 1250, (transient)"
## [1] "Chain 1, iteration 500 of 1250, (sampling)"
## [1] "Chain 1, iteration 750 of 1250, (sampling)"
## [1] "Chain 1, iteration 1000 of 1250, (sampling)"
## [1] "Chain 1, iteration 1250 of 1250, (sampling)"
## [1] "Computing chain 2"
## [1] "Chain 2, iteration 250 of 1250, (transient)"
## [1] "Chain 2, iteration 500 of 1250, (sampling)"
## [1] "Chain 2, iteration 750 of 1250, (sampling)"
## [1] "Chain 2, iteration 1000 of 1250, (sampling)"
## [1] "Chain 2, iteration 1250 of 1250, (sampling)"
```

```
MF = evaluateModelFit(hM=m, predY=preds)
MF$R2
```

```
## [1] 0.3613131
```

We observe that both the explanatory as well as the predictive power is lower than for the model with a spatial random effect.