

HMSC-R 3.0: Getting started with Hmsc-R: spatial models

Gleb Tikhonov Øystein H. Opedal Nerea Abrego Aleksi Lehtikainen
Melinda M. J. de Jonge Jari Oksanen Otso Ovaskainen

28 May 2020

Introduction

The Hierarchical Modelling of Species Communities (HMSC) framework is a statistical framework for analysis of multivariate data, typically from species communities. We assume that the reader has already gone through the vignette “Hmsc-R 3.0: Getting started with Hmsc-R: univariate models” and “Hmsc-R 3.0: Getting started with Hmsc-R: low-dimensional multivariate models”. In the first vignette we shortly discussed how to fit spatially explicit models to the univariate case. Here, we continue to demonstrate how to use HMSC to make spatially explicit models for the multivariate case and for large datasets.

To get Hmsc-R in use, you need to load it and other needed packages.

```
library(Hmsc)
library(MASS)
set.seed(6)
```

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

Generating simulated data

To illustrate how to use spatial models in Hmsc, we generate data for 5 species (**ns**) on 100 sampling units (**n**). We include only one environmental predictor (**x1**) and give the true intercept (**alpha**) and slope (**beta1**) parameters to construct the matrix of the linear predictors. So far, this is similar to what we did in the low-dimensional multivariate case, but note that we now call this linear predictor **Lf** instead of **L**. This indicates that this is the fixed effects part of the linear predictor, i.e. the effect that can be explained by the environmental covariates. In addition to the effect from the environmental covariates, we now add spatially structured residuals which follow the same spatial structure for all species, indicated by **Lr**. To generate these spatial residuals, we first simulate some random x and y coordinates for all sampling units (**xycoords**). We then generate a spatially structured latent predictor (**eta1**) using an exponentially decreasing spatial covariance function where we have set the spatial scale parameter (**alpha**) to 0.35. Next, we set the true slope parameters (**lambda1**) for the species responses to the latent predictor and compute the spatial residuals for each species in **Lr**.

```
n = 100
ns = 5
beta1 = c(-2,-1,0,1,2)
alpha = rep(0,ns)
beta = cbind(alpha,beta1)
x = cbind(rep(1,n),rnorm(n))
Lf = x%*%t(beta)

xycoords = matrix(runif(2*n),ncol=2)
```

```

colnames(xycoords) = c("x-coordinate", "y-coordinate")
rownames(xycoords) = 1:n

sigma.spatial = c(2)
alpha.spatial = c(0.35)
Sigma = sigma.spatial^2*exp(-as.matrix(dist(xycoords))/alpha.spatial)
eta1 = mvrnorm(mu=rep(0,n), Sigma=Sigma)
lambda1 = c(1,2,-2,-1,0)
Lr = eta1%*%t(lambda1)
L = Lf + Lr
y = as.matrix(L + matrix(rnorm(n*ns),ncol=ns))
yprob = 1*((L + matrix(rnorm(n*ns),ncol=ns))>0)
XData = data.frame(x1=x[,2])

```

We can now visualize the species response matrix y as function of the x and y coordinates (Fig. 1). This shows that indeed nearby sampling units have similar responses for species that have a non-zero loading (lambda) to the spatially structured latent variable.

```

rbPal = colorRampPalette(c('cyan', 'red'))
par(mfrow=c(2,3))
Col = rbPal(10)[as.numeric(cut(x[,2],breaks = 10))]
plot(xycoords[,2],xycoords[,1],pch = 20,col = Col,main=paste('x'), asp=1)
for(s in 1:ns){
  Col = rbPal(10)[as.numeric(cut(y[,s],breaks = 10))]
  plot(xycoords[,2],xycoords[,1],pch = 20,col = Col,main=paste('Species',s), asp=1)
}

```

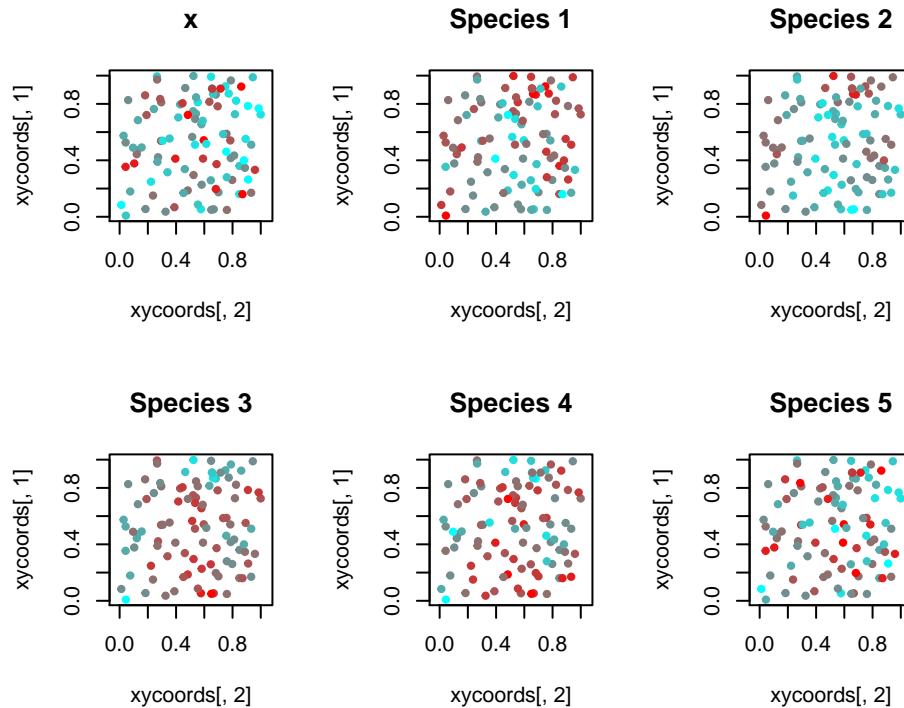


Figure 1: Plots of simulated spatially structured data.

A spatially explicit model in Hmsc

To fit a spatially explicit model with Hmsc, we construct the random effect using the `sData` input argument where we give the coordinates of the sampling units.

```
studyDesign = data.frame(sample = as.factor(1:n))
rL.spatial = HmscRandomLevel(sData = xycoords)
rL.spatial = setPriors(rL.spatial,nfMin=1,nfMax=1) #We limit the model to one latent variables for visu
m.spatial = Hmsc(Y=yprob, XData=XData, XFormula=~x1,
studyDesign=studyDesign, ranLevels=list("sample"=rL.spatial),distr="probit")
```

Model fitting and evaluation of explanatory and predictive power can be done as before. We first set the MCMC sampling parameters.

```
nChains = 2
test.run = FALSE
if (test.run){
  # with this option, the vignette runs fast but results are not reliable
  thin = 1
  samples = 10
  transient = 5
  verbose = 0
} else {
  # with this option, the vignette evaluates slow but it reproduces the results of
  # the .pdf version
  thin = 10
  samples = 1000
  transient = 1000
  verbose = 0
}

m.spatial = sampleMcmc(m.spatial, thin = thin, samples = samples, transient = transient,
  nChains = nChains, verbose = verbose,updater=list(GammaEta=FALSE))
```

```
## Computing chain 1
##Computing chain 2
```

The explanatory and predictive power of the model can now be calculated in the same way as before.

```
#Explanatory power
preds.spatial = computePredictedValues(m.spatial)
MF.spatial = evaluateModelFit(hM=m.spatial, predY=preds.spatial)
MF.spatial
```

```
## $RMSE
## [1] 0.3107951 0.2081848 0.2218778 0.2741564 0.3010669
##
## $AUC
## [1] 0.9335317 0.9886605 0.9850000 0.9518750 0.9504831
##
## $TjurR2
## [1] 0.4604288 0.5815508 0.5689169 0.4527756 0.6197751
```

```
#Predictive power
partition = createPartition(m.spatial, nfolds = 2, column = "sample")
cvpreds.spatial = computePredictedValues(m.spatial, partition=partition,updater=list(GammaEta=FALSE))
```

```
## Cross-validation, fold 1 out of 2
## Computing chain 1
##Computing chain 2
## Cross-validation, fold 2 out of 2
##Computing chain 1
##Computing chain 2

cvMF.spatial = evaluateModelFit(hM=m.spatial, predY=cvpreds.spatial)
cvMF.spatial
```

```
## $RMSE
## [1] 0.3731896 0.3518506 0.3564364 0.3928434 0.3089614
##
## $AUC
## [1] 0.8472222 0.7399008 0.8050000 0.6862500 0.9476651
##
## $TjurR2
## [1] 0.3121820 0.1928824 0.1811694 0.1454237 0.5771773
```

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 random effect is given by the parameter `Alpha[[1]]`. Let's have a look at the MCMC trace plot for this parameter (Fig. 3).

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

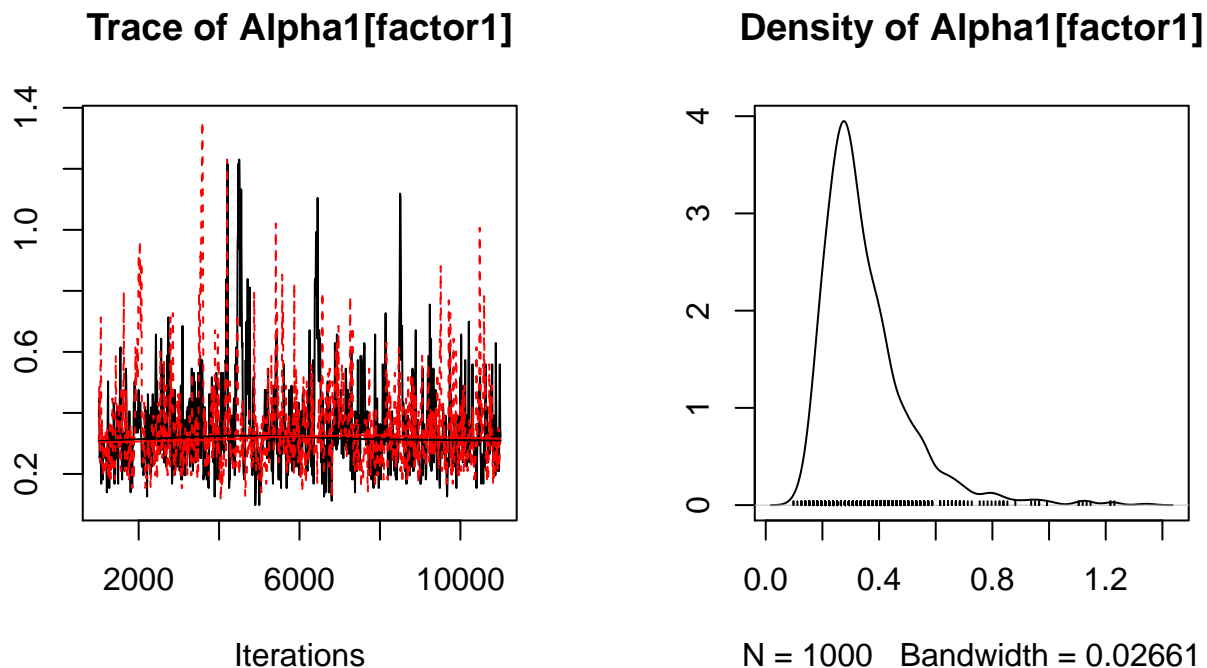


Figure 2: Posterior trace plot of the spatial scale parameter of the spatial model.

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

```
##
```

```
## Iterations = 1010:11000
## Thinning interval = 10
## Number of chains = 2
## Sample size per chain = 1000
##
## 1. Empirical mean and standard deviation for each variable,
##    plus standard error of the mean:
##
##           Mean           SD      Naive SE Time-series SE
##      0.349312      0.155075      0.003468      0.008900
##
## 2. Quantiles for each variable:
##
##   2.5%   25%   50%   75%  97.5%
## 0.1678 0.2517 0.3076 0.4055 0.7691
```

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

```
m = Hmisc(Y=yprob, XData=XData, XFormula=~x1, studyDesign = studyDesign, distr="probit")
m = sampleMcmc(m, thin = thin, samples = samples, transient = transient,
              nChains = nChains, verbose = verbose)
```

```
## Setting updater$GammaEta=FALSE due to absence of random effects included to the model
## Computing chain 1
##Computing chain 2
```

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

```
## $RMSE
## [1] 0.3740062 0.3663433 0.3980664 0.3856937 0.3018576
##
## $AUC
## [1] 0.8462302 0.6661942 0.5793750 0.6743750 0.9504831
##
## $TjurR2
## [1] 0.303236910 0.049953377 0.008418898 0.068500536 0.610422510
```

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

```
## Cross-validation, fold 1 out of 2
## Setting updater$GammaEta=FALSE due to absence of random effects included to the model
## Computing chain 1
##Computing chain 2
## Cross-validation, fold 2 out of 2
## Setting updater$GammaEta=FALSE due to absence of random effects included to the model
##Computing chain 1
##Computing chain 2
```

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

```
## $RMSE
```

```
## [1] 0.3759340 0.3764814 0.4054381 0.3929577 0.3073133
##
## $AUC
## [1] 0.8437500 0.6024096 0.4275000 0.6287500 0.9496779
##
## $TjurR2
## [1] 0.286418211 0.034755606 -0.003307037 0.061799796 0.581567639
```

We observe that both the explanatory and the predictive power is lower than for the model with a spatial random effect for those species that had spatially structured responses.

Spatial models for big datasets

For large datasets, i.e. those with > 1000 sampling units, the standard spatial models as described above may become computationally infeasible. For such datasets, we implemented two alternative approaches to account for the spatial structure in the data: the ‘Nearest Neighbour Gaussian Process (NNGP)’ and the ‘Gaussian Predictive Process (GPP)’. The details of this method are described in Tikhonov et al. (2020).

NNGP models

If we want to fit a NNGP model we have to set `sMethod` to ‘NNGP’ when constructing the random level. Additionally, we can specify how many neighbours we want to use by setting the `nNeighbours` argument. When we do not explicitly set `nNeighbours` when constructing the random level, this parameters is set to 10 as a standard.

```
rL.nngp = HmscRandomLevel(sData = xycoords, sMethod = 'NNGP', nNeighbours = 10)
rL.nngp = setPriors(rL.nngp,nfMin=1,nfMax=1)
```

Running the model and checking the fit is done in the same way as earlier.

```
m.nngp = Hmsc(Y=yprob, XData=XData, XFormula=~x1,
              studyDesign=studyDesign, ranLevels=list("sample"=rL.nngp),distr="probit")
m.nngp = sampleMcmc(m.nngp, thin = thin, samples = samples, transient = transient,
                   nChains = nChains, verbose = verbose, updater=list(GammaEta=FALSE))
```

```
## Computing chain 1
##Computing chain 2
```

```
preds.nngp = computePredictedValues(m.nngp,updater=list(GammaEta=FALSE))
MF.nngp = evaluateModelFit(hM=m.nngp, predY=preds.nngp)
MF.nngp
```

```
## $RMSE
## [1] 0.3099410 0.2006786 0.2318161 0.2795650 0.3009926
##
## $AUC
## [1] 0.9335317 0.9900780 0.9787500 0.9450000 0.9504831
##
## $TjurR2
## [1] 0.4654432 0.6111199 0.5556547 0.4377943 0.6206321
```

```
partition = createPartition(m.nngp, nfolds = 2, column = "sample")
cvpreds.nngp = computePredictedValues(m.nngp, partition=partition,updater=list(GammaEta=FALSE))
```

```
## Cross-validation, fold 1 out of 2
```

```
## Computing chain 1
##Computing chain 2
##Cross-validation, fold 2 out of 2
##Computing chain 1
##Computing chain 2

cvMF.nngp = evaluateModelFit(hM=m.nngp, predY=cvpreds.nngp)
cvMF.nngp
```

```
## $RMSE
## [1] 0.3805924 0.3280701 0.3627340 0.3555614 0.3073361
##
## $AUC
## [1] 0.8358135 0.8802268 0.7587500 0.7975000 0.9472625
##
## $TjurR2
## [1] 0.2831757 0.2161628 0.1960016 0.1789612 0.5881835
```

Let's have a look at the MCMC trace plots for the spatial scale parameter `Alpha[[1]]`.

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

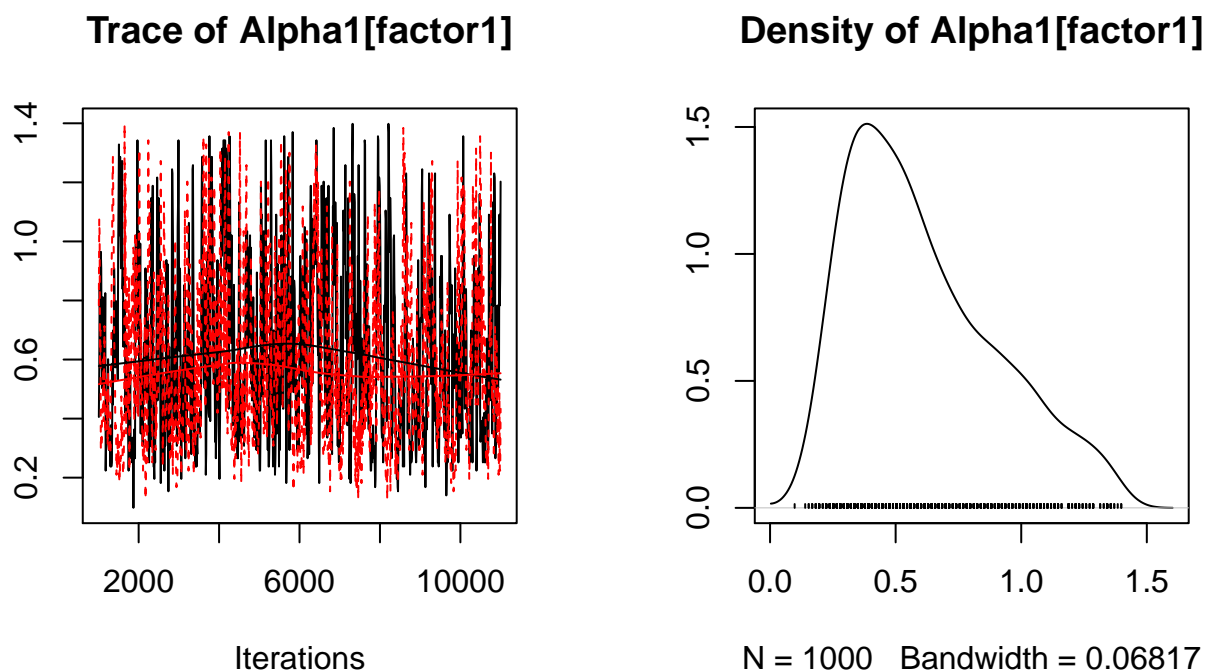


Figure 3: Posterior trace plot of the spatial scale parameter of the nngp spatial model.

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

##
## Iterations = 1010:11000
## Thinning interval = 10
## Number of chains = 2
## Sample size per chain = 1000
##
## 1. Empirical mean and standard deviation for each variable,
```

```
##      plus standard error of the mean:
##
##           Mean           SD      Naive SE Time-series SE
##      0.609467      0.294096      0.006576      0.014718
##
## 2. Quantiles for each variable:
##
##      2.5%      25%      50%      75%      97.5%
## 0.2097 0.3775 0.5453 0.7970 1.2864
```

GPP models

The Gaussian Predictive Process (GPP) assumes that all information on the spatial structure of the data can be summarized at a small number of so called ‘knot’ locations. The locations of these knots have to be specified by the user. To help the user with this, we implemented a function to construct a uniform grid of knots based on the locations of the dataset. This function either need the wanted number of knots along the shortest spatial axis `nKnots` or the wanted distance between knots `KnotDist`. Additionally, the user can specify the maximum distance of a knot to the nearest data point, this ensures that the created grid does contain knots in locations with no datapoints.

```
# Setting the knots
Knots = constructKnots(xycoords, knotDist = 0.2, minKnotDist = 0.4)
```

```
plot(xycoords[,1],xycoords[,2],pch=18, asp=1)
points(Knots[,1],Knots[,2],col='red',pch=18)
```

We can now construct the random level by by setting `sMethod` to ‘GPP’ and supplying the knot locations to `sKnot`.

```
rL.gpp = HmscRandomLevel(sData = xycoords, sMethod = 'GPP', sKnot = Knots)
rL.gpp = setPriors(rL.gpp,nfMin=1,nfMax=1)
m.gpp = Hmsc(Y=yprob, XData=XData, XFormula=~x1,
             studyDesign=studyDesign, ranLevels=list("sample"=rL.gpp),distr="probit")
m.gpp = sampleMcmc(m.gpp, thin = thin, samples = samples, transient = transient,
                  nChains = nChains, verbose = verbose,updater=list(GammaEta=FALSE))
```

```
## Computing chain 1
##Computing chain 2
```

```
preds.gpp = computePredictedValues(m.gpp,updater=list(GammaEta=FALSE))
MF.gpp = evaluateModelFit(hM=m.gpp, predY=preds.gpp)
MF.gpp
```

```
## $RMSE
## [1] 0.3215551 0.1914984 0.2234846 0.2774592 0.3008101
##
## $AUC
## [1] 0.9211310 0.9929128 0.9875000 0.9562500 0.9504831
##
## $TjurR2
## [1] 0.4297787 0.5971137 0.5446531 0.4272606 0.6175804
```

```
cvpreds.gpp = computePredictedValues(m.gpp, partition=partition,updater=list(GammaEta=FALSE))
```

```
## Cross-validation, fold 1 out of 2
##Computing chain 1
```

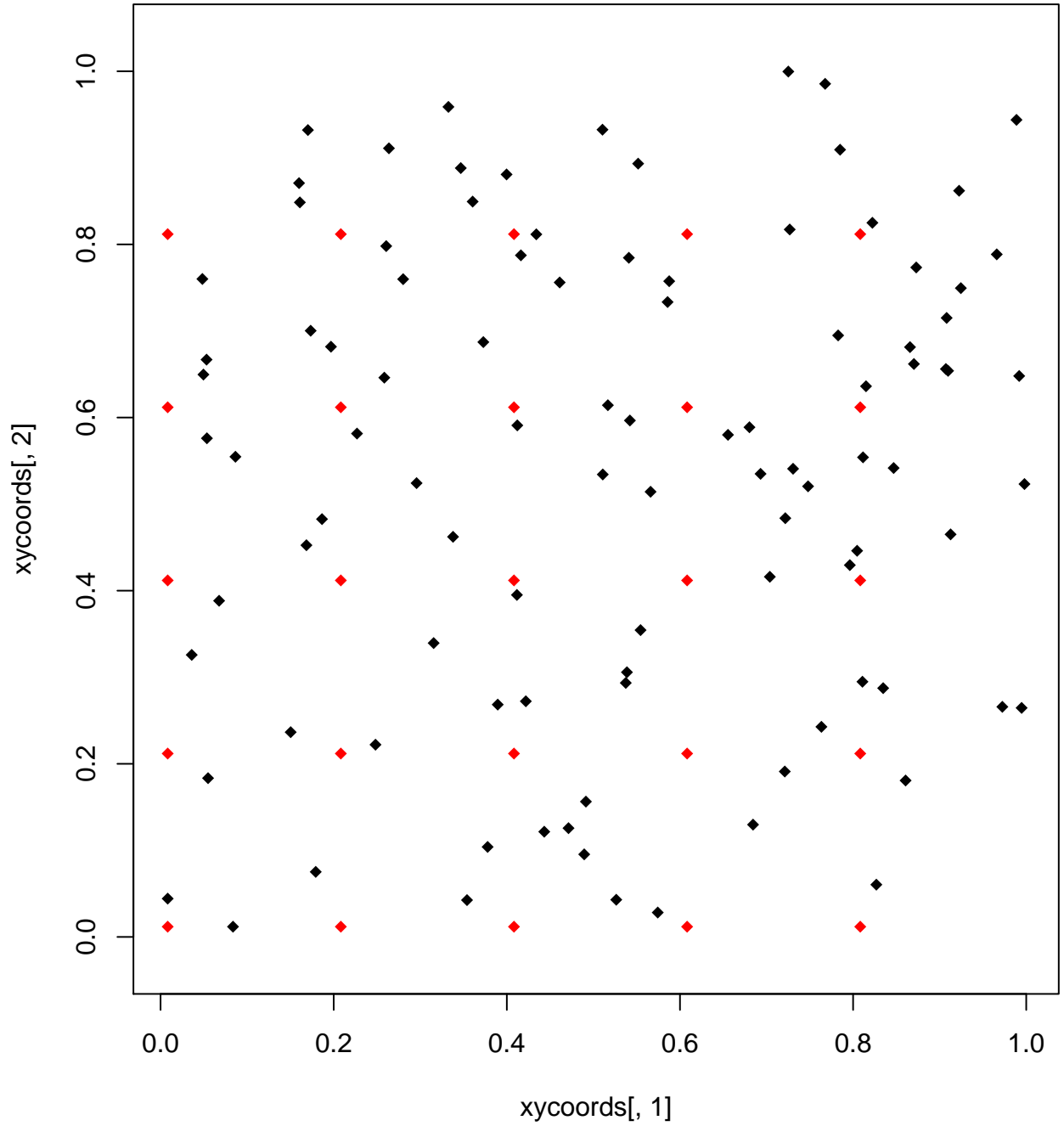



Figure 4: Locations of the created knots in red together with the locations of the plots in black.

```
## Computing chain 2
## Cross-validation, fold 2 out of 2
## Computing chain 1
## Computing chain 2

cvMF.gpp = evaluateModelFit(hM=m.gpp, predY=cvpreds.gpp)
cvMF.gpp

## $RMSE
## [1] 0.3891450 0.3533414 0.3873238 0.3798896 0.3078829
##
## $AUC
## [1] 0.8229167 0.8029766 0.6906250 0.7306250 0.9488728
##
## $TjurR2
## [1] 0.25696389 0.10581360 0.07826134 0.09500601 0.58461732

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

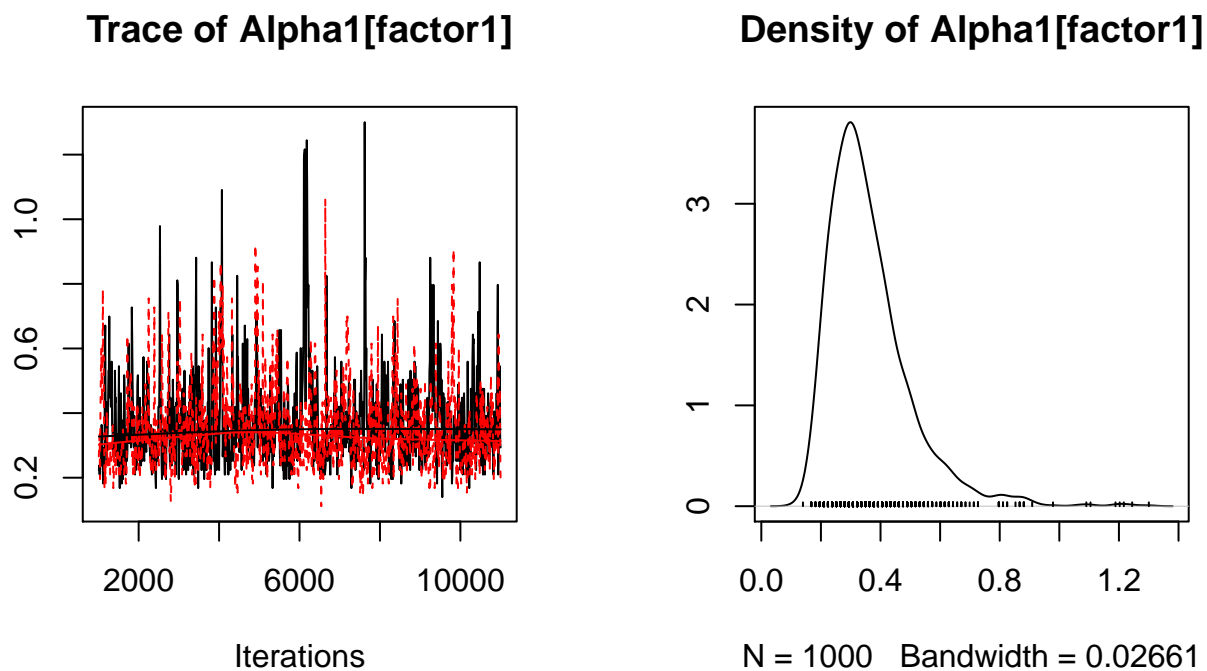


Figure 5: Posterior trace plot of the spatial scale parameter of the ggp spatial model.

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

##
## Iterations = 1010:11000
## Thinning interval = 10
## Number of chains = 2
## Sample size per chain = 1000
##
## 1. Empirical mean and standard deviation for each variable,
##    plus standard error of the mean:
##
##           Mean           SD      Naive SE Time-series SE
```

```
##      0.361456      0.139139      0.003111      0.006803
##
## 2. Quantiles for each variable:
##
##   2.5%   25%   50%   75%  97.5%
## 0.1818 0.2657 0.3356 0.4195 0.7131
```

For this small dataset of 200 sampling units, you will notice that the sampling times for the three spatial models are very similar. However, for larger datasets, the sampling times of the NNGP and GPP method are much shorter than for the standard spatial model.

As a last remark, you may have noticed that we set the `updater` parameter in `sampleMCMC` when running the NNGP and GPP method. With this parameter we specify which MCMC updaters we want to include during the MCMC sampling. The standard setting is to use all available updaters. However, the GammaEta updater is currently not available for these methods. This is not a problem however because this is an optional updater which helps reach convergence in less iterations in some situations.

References

Tikhonov, G., L. Duan, N. Abrego, G. Newell, M. White, D. Dunson, and O. Ovaskainen. 2020. “Computationally Efficient Joint Species Distribution Modeling of Big Spatial Data.” *Ecology* 101: e02929. doi:10.1002/ecy.2929.