

# Hierarchical Bayesian species distribution models with the **hSDM** R Package



*Adansonia grandidieri* Baill. next to Andavadoaka village (southwest Madagascar).

Ghislain Vieilledent<sup>\*,1</sup>

[\*] **Corresponding author:** \E-mail: [ghislain.vieilledent@cirad.fr](mailto:ghislain.vieilledent@cirad.fr) \Phone: +33.(0)4.67.59.37.51  
\Fax: +33.(0)4.67.59.39.09

[1] **Cirad** – UPR BSEF, F-34398 Montpellier, France

## **Abstract**

Work in progress...

*Keywords:* R

# 1 Introduction

## 2 Species distribution models

### 2.1 Binomial model

#### 2.1.1 Mathematical formulation

Let's consider a random variable  $y_{i(jt)}$  (abbreviated  $y_i$ ) representing the total number of presences of a species after several visits  $v_{i(jt)}$  (abbreviated  $v_i$ ) at a particular site  $j$  at time  $t$  ( $t$  can indicate a day, a year, etc. thus allowing several visits at time  $t$ ). Random variable  $y_i$  can take values from 0 to  $v_i$  and can be assumed to follow a Binomial distribution having parameters  $v_i$  and  $\theta_{i(jt)}$  (Eq. 2).  $\theta_{i(jt)}$  (abbreviated  $\theta_i$ ) can be interpreted as the probability of presence of the species at site  $j$  and time  $t$ . Using a logit link function,  $\theta_i$  can be expressed as a linear model combining explicative variables  $X_{i(jt)}$  and parameters  $\beta$  (Eq. 2).

$$(1) \quad \begin{aligned} y_i &\sim \text{Binomial}(v_i, \theta_i) \\ \text{logit}(\theta_i) &= X_{i(jt)}\beta \end{aligned}$$

Using this statistical model, we aim at representing a “suitability process”. Given environmental variables  $X_{i(jt)}$ , how much is habitat at site  $j$  and at time  $t$  suitable for the species under consideration? Parameters  $\beta$  indicate how much each environmental variable contributes to the suitability process. Like every other function in the **hSDM** R package, function `hSDM.binomial()` estimates the parameters  $\beta$  of such a model in a Bayesian framework. Parameter inference is done using a Gibbs sampler including a Metropolis algorithm. The Gibbs sampler is coded in the C language to optimize computation efficiency.

#### 2.1.2 Data generation

To explore the characteristics of the `hSDM.binomial()` function, we can generate a virtual data-set on the basis of the Binomial model described above (Eq. 2). In the most general case (presence/absence data or presence/pseudo-absence data), a site  $j$  is visited once at time  $t$  and  $v_i = 1$ . Thus, the random variable  $y_i$  follows a Bernoulli distribution having parameters  $\theta_i$  and habitat characteristics  $X_{i(jt)}$  are fixed for site  $j$ . We will generate a virtual data-set in this particular case. For data generation, we can import virtual altitudinal data in R. Altitude will be used as an explicative variable to determine habitat suitability, i.e. the probability of presence of a virtual species. Altitudinal data are available on the website of the **hSDM** R package hosted on Sourceforge (<http://hsdm.sourceforge.net/altitude.csv>).

These data can be transformed into a raster using function `rasterFromXYZ()` from the **raster** package. The raster has 2500 cells (50 columns and 50 rows) and the altitude

is comprised roughly between 100 and 600 m (Fig. 7). For linear models, explicative variables are usually centered and scaled to facilitate inference and interpretation of model parameters.

```
# Import altitudinal data
library(raster)
fname <- "http://hsdm.sourceforge.net/altitude.csv"
alt.df <- read.csv(fname,header=TRUE)
alt.orig <- rasterFromXYZ(alt.df)
plot(alt.orig)
# Center and scale altitudinal data
alt <- scale(alt.orig,center=TRUE,scale=TRUE)
plot(alt)
```

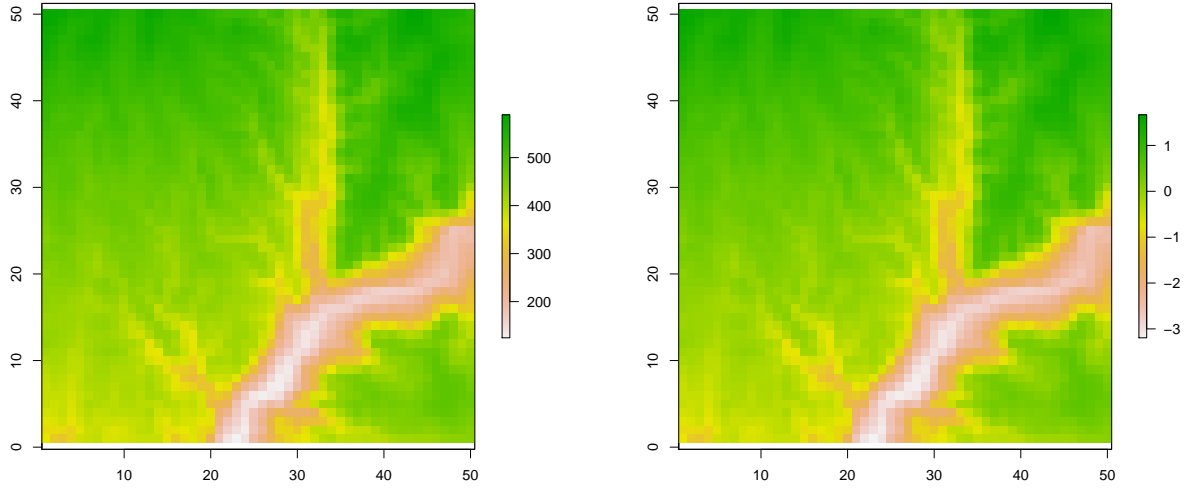


Figure 1: **Altitudinal data.** Original values (in m) on the left. Centered and scaled values on the right.

A quadratic effect of the altitude (variable denoted  $x$ ) is used to compute the probability of presence of the species (Eq. 2).

$$(2) \quad y_i \sim \text{Bernoulli}(\theta_i)$$

$$\text{logit}(\theta_i) = \beta_0 + \beta_1 x_i + \beta_2 x_i^2$$

We fix the parameters to  $\beta_0 = 1$ ,  $\beta_1 = 2$  and  $\beta_2 = -4$ . The species has a higher probability of presence at intermediate altitudinal values (Fig. 2).

```

# Load hSDM library
library(hSDM)
# Target parameters
beta.target <- matrix(c(1,2,-4),ncol=1)
# Matrix of covariates (including the intercept)
X <- cbind(rep(1,ncell(alt)),values(alt),(values(alt))^2)
# Probability of presence as a quadratic function of altitude
logit.theta <- X %*% beta.target
theta <- inv.logit(logit.theta)
# Transform the probability of presence into a raster
theta <- rasterFromXYZ(cbind(coordinates(alt),theta))
# Plot the probability of presence
plot(theta,col=rev(heat.colors(255)))

```

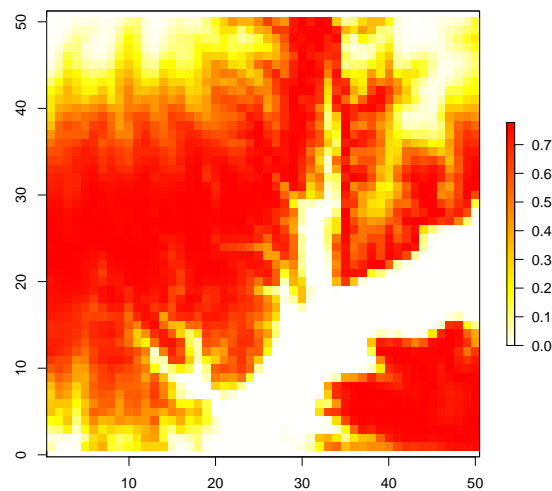


Figure 2: **Probability of presence.**

We can assume a number  $n$  of points in the landscape where we have been able to observe or not the presence of the species. We can simulate the presence or absence of the species at these  $n$  points given our model (Fig. 3).

```

# Load dismo library
library(dismo) # For randomPoints() function
# Number of observation points
nobspt <- 200
# Set seed for repeatability
seed <- 1234

```

```

# Sample the observations in the landscape
set.seed(seed)
obs <- randomPoints(alt,nobspt)
# Extract altitude data for observations
alt.obs <- extract(alt,obs)
# Compute theta for these observations
X.obs <- cbind(rep(1,nobspt),alt.obs,alt.obs^2)
logit.theta.obs <- X.obs %*% beta.target
theta.obs <- inv.logit(logit.theta.obs)
# Simulate observations
trials <- rep(1,nobspt)
set.seed(seed)
Y <- rbinom(nobspt,trials,theta.obs)
# Group explicative and response variables in a data-frame
data.obs.df <- data.frame(Y,trials=trials,alt=alt.obs)
# Transform observations in a spatial object
library(sp)
data.obs <- SpatialPointsDataFrame(coords=coordinates(obs),data=data.obs.df)
# Plot observations
plot(alt.orig)
points(data.obs[data.obs$Y==1,],pch=16)
points(data.obs[data.obs$Y==0,],pch=1)

```

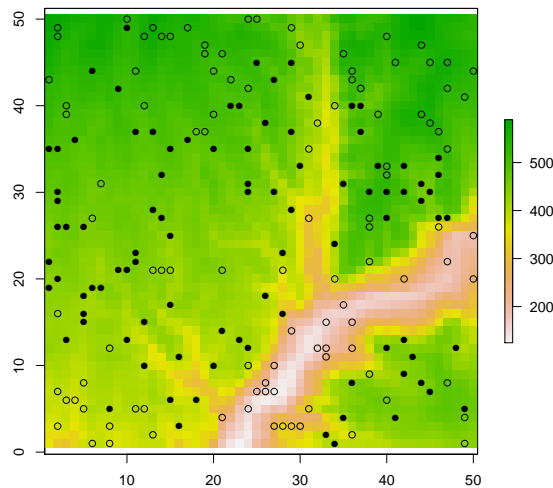


Figure 3: **Observation points.** Presences (full circles) and absences (empty circles) are localized on the altitude map (in m).

### 2.1.3 Parameter inference using the `hSDM.binomial()` function

The `hSDM.binomial()` function performs a Binomial logistic regression in a Bayesian framework. Before using this function we need to prepare a bit the data for parameter inference and prediction. For parameter inference, we add the quadratic term for altitude in the data frame associated to observations.

```
data.obs$alt2 <- (data.obs$alt)^2
```

We want to have predictions on the whole landscape, not only at observation points. To directly obtain these predictions, we can create a data frame including altitudinal data on the whole landscape. This data frame will be used for the `suitability.pred` argument. The data frame for predictions must include the same column names as those used in the formula for the `suitability` argument (i.e. “alt” and “alt2” in our example).

```
data.pred <- data.frame(alt=values(alt),alt2=(values(alt))^2)
```

We can now call the `hSDM.binomial()` function. Setting parameter `save.p` to 1, we can save in memory the MCMC values for predictions. These values can be used to compute several statistics for each predictions (mean, median, 95% quantiles). For example, mean and 95% quantiles are useful to estimate the uncertainty around the mean predictions.

```
mod.hSDM.binomial <- hSDM.binomial(presences=data.obs$Y,  
                                   trials=data.obs$trials,  
                                   suitability=~alt+alt2,  
                                   data=data.obs,  
                                   suitability.pred=data.pred,  
                                   burnin=1000, mcmc=1000, thin=1,  
                                   beta.start=0,  
                                   mubeta=0, Vbeta=1.0E6,  
                                   seed=1234, verbose=1, save.p=1)
```

### 2.1.4 Analysis of the results

The `hSDM.binomial()` function returns an MCMC (Markov chain Monte Carlo) for each parameter of the model and also for the model deviance. To obtain parameter estimates, MCMC values can be summarized through a call to the `summary()` function from the **coda** package. We can check that the values of the target parameters  $\beta_0 = 1$ ,  $\beta_1 = 2$  and  $\beta_2 = -4$  are within the 95% confidence interval of the parameter estimates.

```
summary(mod.hSDM.binomial$mcmc)
```

```
##
## Iterations = 1001:2000
## Thinning interval = 1
## Number of chains = 1
## 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
## beta.(Intercept)  1.18 0.231  0.00732      0.0215
## beta.alt          1.57 0.556  0.01758      0.0821
## beta.alt2         -3.89 0.694  0.02194      0.1116
## Deviance          196.27 2.430  0.07685      0.2225
##
## 2. Quantiles for each variable:
##
##              2.5%    25%    50%    75%    97.5%
## beta.(Intercept)  0.735    1.02    1.16    1.34    1.62
## beta.alt          0.468    1.16    1.57    1.97    2.64
## beta.alt2         -5.200   -4.37   -3.88   -3.37   -2.41
## Deviance          193.428 194.50 195.67 197.53 202.85
```

Parameters estimates can be compared to results obtained with the `glm()` function.

```
##= glm results for comparison
mod.glm <- glm(cbind(Y,trials-Y)~alt+alt2,family="binomial",data=data.obs)
summary(mod.glm)

##
## Call:
## glm(formula = cbind(Y, trials - Y) ~ alt + alt2, family = "binomial",
##      data = data.obs)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -1.766   -0.797    0.000    0.799    3.112
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)    1.186     0.241    4.92 8.6e-07 ***
## alt           1.475     0.538    2.74  0.0061 **
## alt2          -3.736     0.735   -5.08 3.7e-07 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```



```
##
## (Dispersion parameter for binomial family taken to be 1)
##
##      Null deviance: 276.54  on 199  degrees of freedom
## Residual deviance: 193.15  on 197  degrees of freedom
## AIC: 199.1
##
## Number of Fisher Scoring iterations: 8
```

MCMC can also be graphically summarized with a call to the `plot.mcmc()` function, also in the **coda** package. MCMC are plotted with a trace of the sampled output and a density estimate for each variable in the chain (Fig. 4). This plot can be used to visually check that the chains have converged.

```
plot(mod.hSDM.binomial$mcmc)
```

The `hSDM.binomial()` function also returns two other objects. The first one, `prob.p.latent`, is the predictive posterior mean of the latent variable  $\theta$  (the probability of presence) for each observation.

```
str(mod.hSDM.binomial$prob.p.latent)

##  num [1:200] 0.17046 0.44498 0.6823 0.00542 0.05014 ...

summary(mod.hSDM.binomial$prob.p.latent)

##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##      0.000   0.184   0.572   0.465   0.731   0.789
```

The second one, `prob.p.pred` is the set of sampled values from the predictive posterior (if parameter `save.p` is set to 1) or the predictive posterior mean (if `save.p` is set to 0) for each prediction. In our example, `save.p` is set to 1 and `prob.p.pred` is an `mcmc` object. Values in `prob.p.pred` can be used to plot the predicted probability of presence on the whole landscape and the uncertainty associated to the predictions (Fig 5).

```
# Create a raster for predictions
theta.pred.mean <- raster(theta)
# Create rasters for uncertainty
theta.pred.2.5 <- theta.pred.97.5 <- raster(theta)
# Attribute predicted values to raster cells
theta.pred.mean[] <- apply(mod.hSDM.binomial$prob.p.pred,2,mean)
theta.pred.2.5[] <- apply(mod.hSDM.binomial$prob.p.pred,2,quantile,0.025)
theta.pred.97.5[] <- apply(mod.hSDM.binomial$prob.p.pred,2,quantile,0.975)
# Plot the predicted probability of presence and uncertainty
```

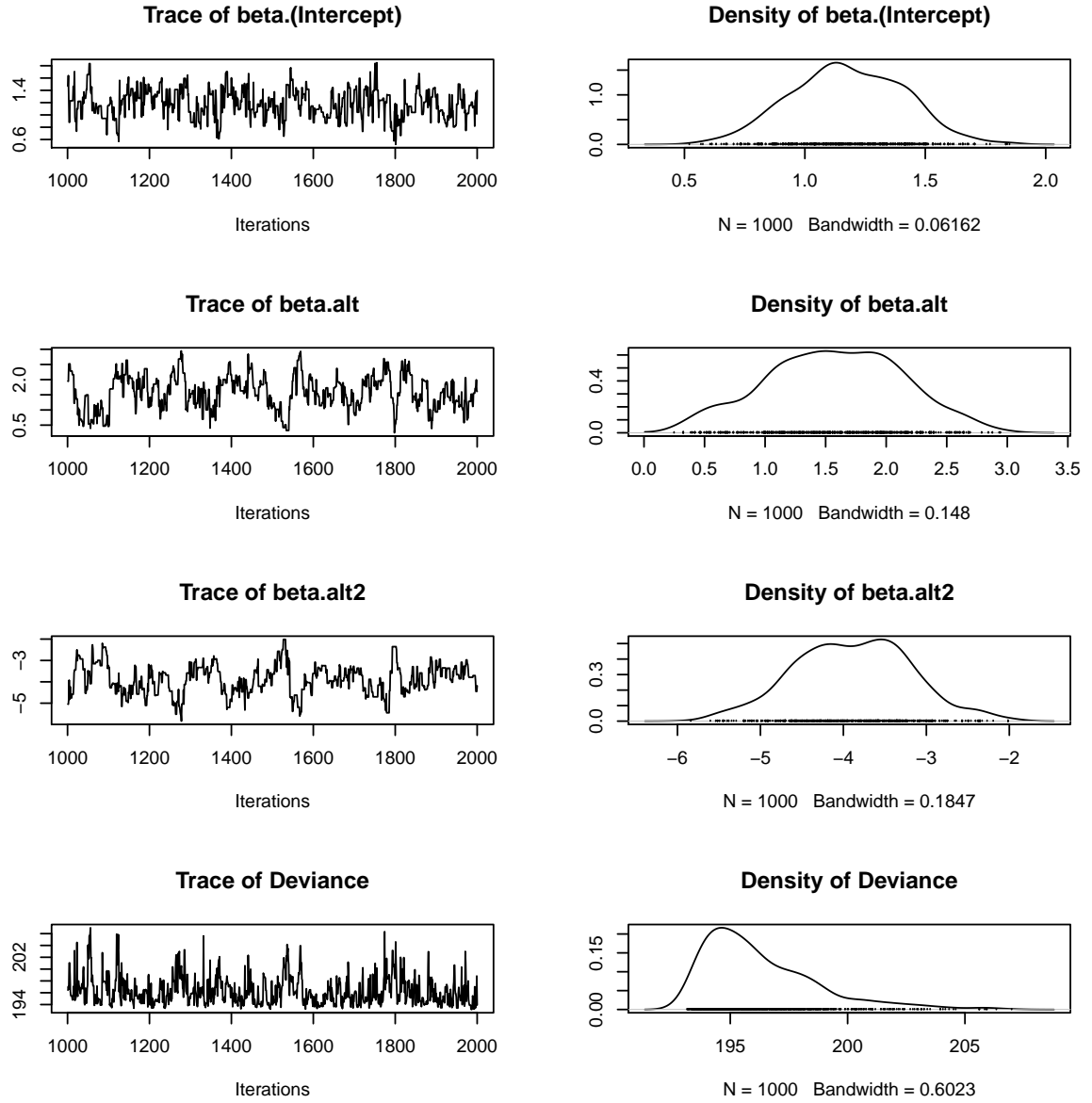


Figure 4: Trace and density estimate for each variable of the MCMC.

```
plot(theta.pred.mean,col=rev(heat.colors(255)))
plot(theta.pred.2.5,col=rev(heat.colors(255)))
plot(theta.pred.97.5,col=rev(heat.colors(255)))
```

In our example, we can compare the predictions to the initial probability of presence computed from our model to check that our predictions are correct (Fig. 6).

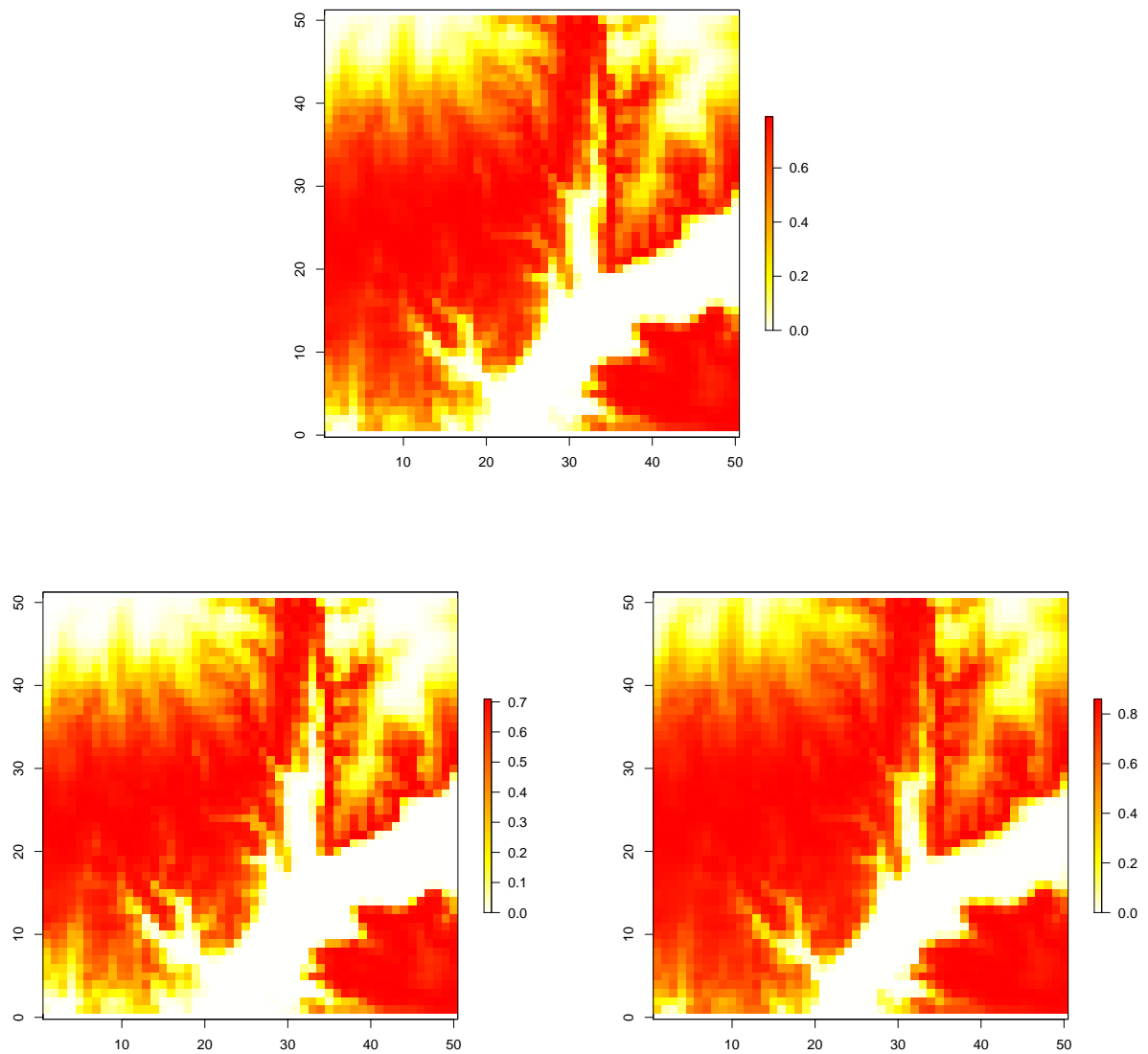


Figure 5: **Predicted probability of presence and uncertainty of predictions.** Mean probability of presence (top), predictions at 2.5% quantile (bottom left) and 97.5% quantile (bottom right) can be plotted from the `mcmc` object `plot.p.pred` returned by function `hSDM.binomial()`.

```
# Comparing predictions to initial values
plot(theta[], theta.pred.mean[], cex.lab=1.4)
abline(a=0, b=1, col="red", lwd=2)
```

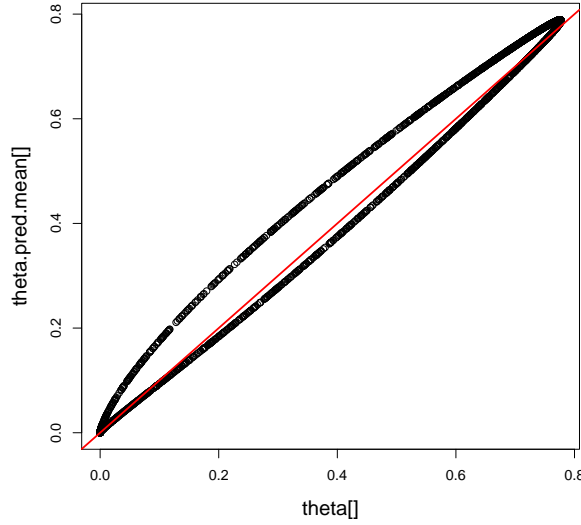


Figure 6: **Predicted vs. initial probabilities of presence.** Initial probabilities of presence are computed from the Binomial logistic regression model with fixed parameters.

## 2.2 Site-occupancy model

### 2.2.1 Mathematical formulation

Occurrence of a species is typically not observed perfectly. Species traits, survey-specific conditions and site-specific characteristics may influence species detection probability which is often  $< 1$  (Chen *et al.*, 2013). Thus, observations might include false absences. For example, the habitat can be suitable and the species present but individuals have not been seen during the census. Or the habitat can be suitable but the species has not dispersed yet to the site (typical example for plant species, see Latimer *et al.* (2006)) or was not present on the site at the moment of the observation (typical example for animal species such as birds, see Kéry *et al.* (2005)). Treating observed occurrence and species distributions as the true occurrence and distribution, failing to make amendments for imperfect detection, may lead to problems in species distribution studies, habitat models and biodiversity management (Kéry & Schmidt, 2008; Lahoz-Monfort *et al.*, 2014; Latimer *et al.*, 2006).

New classes of models, called site-occupancy models, were developed to solve the problems created by imperfect detection (MacKenzie *et al.*, 2002). These models combine two processes, an ecological process which describes habitat suitability and an observation process which takes into account imperfect detection. Site-occupancy models use information from repeated observations at each site to estimate detectability. Detectability may vary with site characteristics (e.g., habitat variables) or survey characteristics (e.g., weather conditions), whereas occupancy relates only to site characteristics.

Let's consider the random variable  $z_{i(jt)}$  (abbreviated  $z_i$ ) describing habitat suitability for observation  $i$  at site  $j$  and at time  $t$ . The random variable  $z_i$  can take value 1 or

0 depending on the fact that the habitat at site  $j$  and at time  $t$  is suitable ( $z_i = 1$ ) or not ( $z_i = 0$ ). Habitat at site  $j$  and time  $t$  is described by environmental variables  $X_{i(jt)}$ . Random variable  $z_i$  can be assumed to follow a Bernoulli distribution of parameter  $\theta_i$  (Eq. 3). In this case,  $\theta_i$  is the probability that the habitat is suitable. Let's consider also the random variable  $y_{i(jt)}$  (abbreviated  $y_i$ ) representing the total number of presences of a species after several visits  $v_{i(jt)}$  (abbreviated  $v_i$ ) at a particular site  $j$  at time  $t$ . Again,  $t$  can indicate a day, a year, etc. thus allowing several visits at time  $t$ . The species is observed ( $y_i \geq 1$ ) only if the habitat is suitable ( $z = 1$ ). The species is unobserved ( $y_i = 0$ ) if the habitat is not suitable ( $z_i = 0$ ) or if the habitat is suitable ( $z_i = 1$ ) but the detection probability  $\delta_{i(jt)}$  (abbreviated  $\delta_i$ ) is inferior to 1. Thus, the total number of presences  $y_i$  is assumed to follow a Binomial distribution with parameters  $z_i\delta_i$  and  $v_i$ . Using a logit link function,  $\delta_i$  can be expressed as a linear model combining explicative variables  $W_{i(jt)}$  and parameters  $\gamma$  (Eq. 3). Typically, explicative variables  $W_{i(jt)}$  are site characteristics (e.g., habitat variables) or survey characteristics (e.g., weather conditions). The function `hSDM.siteocc()` estimates the parameters  $\beta$  and  $\gamma$  of such a model.

**Ecological process:**

$$z_i \sim \mathcal{Bernoulli}(\theta_i)$$

$$\text{logit}(\theta_i) = X_{i(jt)}\beta$$

(3)

**Observation process:**

$$y_i \sim \mathcal{Binomial}(z_i\delta_i, v_i)$$

$$\text{logit}(\delta_i) = W_{i(jt)}\gamma$$

### 2.2.2 Data generation

To explore the characteristics of the `hSDM.siteocc()` function, we can generate a new virtual data-set on the basis of the site-occupancy model described above (Eq. 3). In most general cases, protocol experiments would include several visits with varying survey conditions (e.g. weather conditions) to several sites with fixed sites characteristics (e.g. habitat variables). We will generate virtual data following this protocole and using the altitudinal data used in the previous example for the Binomial model (Sec. 2.1).

We draw at random the number of visits at each observation point used in the previous example (see Fig. 3 of Sec. 2.1).

```
# Number of visits associated to each observation point
set.seed(seed)
v <- rpois(nobspt,2)
v[v==0] <- 1
```

We fix the survey conditions for each visit depending on two explicative variables  $w_1$  and  $w_2$  which will explain the observability of the species. We also fix the intercept and

the effects of these two variables:  $\gamma_0 = 0.2$ ,  $\gamma_1 = 0.5$  and  $\gamma_2 = 0.5$  for determining the detection probability.

```
# Explicative variables for observation process
nobs <- sum(v)
W <- cbind(rep(1,nobs))
set.seed(seed)
w1 <- rnorm(n=nobs,0,1)
set.seed(2*seed)
w2 <- rnorm(n=nobs,0,1)
W <- cbind(rep(1,nobs),w1,w2)
# Target parameters for observation process
gamma.target <- matrix(c(0.2,0.5,0.5),ncol=1)
```

Using covariates and parameters for the two processes, we compute the probability that the habitat is suitable ( $\theta_i$ ) and the species detection probability ( $\delta_i$ ). We also draw the random variables  $z_i$  and  $y_i$  and construct the observation data-set.

```
# Ecological process (suitability)
logit.theta <- X %*% beta.target
theta <- inv.logit(logit.theta)
set.seed(seed)
z.obspt <- rbinom(nobspt,1,theta)
Z <- rep(z.obspt,each=v)

# Observation process (detectability)
logit.delta <- W %*% gamma.target
delta <- inv.logit(logit.delta)
set.seed(seed)
Y <- rbinom(nobs,1,delta*Z)

# Data-set
X1 <- rep(X[,2],each=v)
X2 <- rep(X[,3],each=v)
trials <- rep(1,nobs)
data.obs <- data.frame(Y,trials,X1,X2,W1,W2)

## Error: objet 'W1' introuvable
```

## 2.3 Binomial iCAR model

## 2.4 Site-occupancy iCAR model

### 3 Acknowledgements

## References

- Chen G, Kéry M, Plattner M, Ma K, Gardner B (2013) Imperfect detection is the rule rather than the exception in plant distribution studies. *Journal of Ecology*, **101**, 183–191.
- Kéry M, Royle JA, Schmid H (2005) Modeling avian abundance from replicated counts using binomial mixture models. *Ecological applications*, **15**, 1450–1461.
- Kéry M, Schmidt BR (2008) Imperfect detection and its consequences for monitoring for conservation. *Community Ecology*, **9**, 207–216.
- Lahoz-Monfort JJ, Guillera-Arroita G, Wintle BA (2014) Imperfect detection impacts the performance of species distribution models. *Global Ecology and Biogeography*, **23**, 504–515. doi:10.1111/geb.12138. URL <http://dx.doi.org/10.1111/geb.12138>.
- Latimer AM, Wu SS, Gelfand AE, Silander JA (2006) Building statistical models to analyze species distributions. *Ecological Applications*, **16**, 33–50.
- MacKenzie DI, Nichols JD, Lachman GB, Droege S, Andrew Royle J, Langtimm CA (2002) Estimating site occupancy rates when detection probabilities are less than one. *Ecology*, **83**, 2248–2255.