

# Accounting for spatial regional variability in modelling and forecasting deforestation – the fate of Madagascar’s forests

Ghislain VIEILLEDENT<sup>1,2,3,4,★</sup> and Frédéric ACHARD<sup>1</sup>

[1] **European Commission** – JRC, Bio-economy Unit, I-21027 Ispra (VA), ITALY

[2] **CIRAD** – UPR Forêts et Sociétés, F-34398 Montpellier, FRANCE

[3] **CIRAD** – UPR AMAP, F-34398 Montpellier, FRANCE

[4] **Univ Montpellier** – AMAP, CIRAD, Montpellier, FRANCE

[★] **Corresponding author:** \E-mail: ghislain.vieilledent@cirad.fr \Phone: +33 4 67 61 49 09

# Abstract

Deforestation models are useful tools in landscape ecology. They can be used to identify the main drivers of deforestation and estimate their relative contribution. When spatially explicit, models can also be used to predict the location of future deforestation. Deforestation forecasts can be used to estimate associated CO<sub>2</sub> emissions responsible of climate change, prioritize areas for conservation and identify refuge areas for biodiversity. Most of spatial deforestation models includes landscape variables such as the distance to forest edge, the distance to nearest road or the presence of protected areas. Such variables commonly explain a small part of the deforestation process and a large spatial variability remains unexplained by the model.

In the present study, we show how using an intrinsic conditional autoregressive (iCAR) model in a hierarchical Bayesian approach can help structure the residual spatial variability in the deforestation process and obtain more realistic predictions of the deforestation (Dormann et al. 2007). We take Madagascar as a case study considering deforestation on the period 1990-2010 and forecasting deforestation on the period 2010-2050. We demonstrate that accounting for spatial autocorrelation increases the percentage of explained deviance of 21 points for the deforestation model in Madagascar. We also illustrate the use of the newly developed **forestatrisk** Python module to rapidly estimate the numerous parameters of a deforestation model including an iCAR process and efficiently forecast deforestation on large geographical areas at high spatial resolution.

We advocate the use of such models to obtain more accurate land-use change predictions. Such an approach could be used to estimate better the impact of future deforestation in the global carbon cycle and define more efficient strategies for biodiversity conservation in tropical countries.

# 1 Introduction

## 2 Materials and Methods

### 2.1 Data

We used historical deforestation maps for Madagascar at 30m resolution for three time-periods: 1990-2000, 2000-2010, and 2010-2017 (Vieilledent *et al.*, 2018b). We tried to model the observed spatial deforestation process on the period 2000-2010 at the national level. Period 1990-2000 was used to compute the distance to past deforestation for each forest pixel in 2000. Period 2010-2017 was used to compare model forecasts with observations.

To explain the observed spatial deforestation on the period 2000-2010, we considered various spatial explanatory variables describing: topography (altitude and slope), accessibility (distances to nearest road, town and river), forest landscape (distance to forest edge), deforestation history (distance to past deforestation) and land-tenure variables (protected area system). Characteristics of each variables are summarized in Tab. 1.

Altitude (in m) and slope (in degree) at 90m resolution were obtained from the SRTM Digital Elevation Database v4.1 (<http://srtm.csi.cgiar.org/>). Distances (in m) to nearest road, town and river at 150m resolution were derived from the OpenStreetMap (OSM) project for Madagascar (<http://www.geofabrik.de/>). To obtain the road network in Madagascar, we considered the categories “motorway”, “trunk”, “primary”, “secondary” and “tertiary” for the “highway” key in OSM. To obtain the network of populated places in Madagascar (that we simply call “towns” in the present study), we considered the categories “city”, “town” and “village” for the “place” key in OSM. To obtain the river network in Madagascar, we considered the categories “river” and “canal” for the “waterway” key in OSM. For a more detailed description of each category, see the OSM wiki page (<https://wiki.openstreetmap.org/wiki/Tags>). Distance to forest edge was computed at 30m resolution from the forest cover map in 2000. Distance to past deforestation was computed at 30m resolution from the 1990-2000 forest cover change map. For the protected area system, we used the 20/12/2010 version of the SAPM “Système des Aires Protégées à Madagascar” (<http://rebioma.net/>) and considered both Protected Areas (created before 2003) and New Protected Areas in the SAPM terminology. Polygons representing protected areas were rasterized at 30m resolution. In total, we obtained 8 spatial variables that could explained the location of the deforestation.

### 2.2 Models

We compared two deforestation models. The first model described in Eq. (1) is a simple logistic regression model, a special case of generalized linear model (GLM) for binary data. This model is denoted “glm” in subsequent sections and results. We considered the random variable  $y_i$  which takes value 1 if the forest pixel  $i$  is deforested on the period 2000-2010 and 0 if it is not. We assumed that  $y_i$  follows a Bernoulli distribution of parameter  $\theta_i$ . In our model,  $\theta_i$  represents the spatial probability of deforestation for pixel  $i$ . We assumed that  $\theta_i$  is linked, through a logit function, to a linear combination of the explanatory variables

$X_i\beta$ , where  $X_i$  is the vector of explanatory variables for pixel  $i$  and  $\beta$  is the vector of model parameters to be estimated.

$$\begin{aligned} y_i &\sim \text{Bernoulli}(\theta_i) \\ \text{logit}(\theta_i) &= X_i\beta \end{aligned} \tag{1}$$

The second model described in Eq. (2) includes additional random effects  $\rho_j$  for each spatial cell  $j$  of a  $10 \times 10$  km grid covering Madagascar. This grid resolution was chosen in order to have a reasonable balance between a good representation of the spatial variability of the deforestation process and the number of parameters. We assumed that random effects were spatially autocorrelated through an intrinsic conditional autoregressive (iCAR) model (Besag *et al.*, 1991; Banerjee *et al.*, 2014). This model is denoted “icar” in subsequent sections and results. In a iCAR model, the random effect  $\rho_j$  associated to cell  $j$  depends on the values of the random effects  $\rho_{j'}$  associated to neighbouring cells  $j'$ . In our case, the neighbouring cells are cells connected to the target cell  $j$  through a common border or corner (cells defined by the “king move” in chess). The number of neighbouring cells for cell  $j$ , which might vary, is denoted  $n_j$ .

$$\begin{aligned} y_i &\sim \text{Bernoulli}(\theta_i) \\ \text{logit}(\theta_i) &= X_i\beta + \rho_j \\ \rho_j &\sim \text{Normal}(\sum_{j'} \rho_{j'}/n_j, V_\rho/n_j) \end{aligned} \tag{2}$$

The first model can be viewed as a “process-based” model for which variables are selected on an *a priori* knowledge of the deforestation process. For example, we assumed that the risk of deforestation decreases with the distance to road and forest edge, and is lower in protected areas. The second model can be viewed as a model combining a “process-based” part and a “pattern-oriented” part. Additional spatial random effects  $\rho_j$  account for unmeasured or unmeasurable variables that explain a part of the residual spatial variation in the deforestation process (the residual spatial “pattern”) that is not explained by the fixed environmental variables ( $X_i$ ). While the first model has only 9 parameters to be estimated (one intercept parameter plus 8 slope parameters for the explanatory variables), the second model has 6,266 parameters to be estimated, including the 6,257 spatial random effects for the  $10 \times 10$  km cells covering whole Madagascar (for which lands cover 587 000 km<sup>2</sup>)

We used a random sample of 20,000 forest pixels in 2000 to fit the two models. The sample was stratified between 10,000 deforested pixels in 2000-2010 and 10,000 non-deforested pixels. A balanced sample between deforested and non-deforested pixels is preferable in our case (Dezécache *et al.*, 2017). First, deforestation events are rare ( $\sim 1$  %/yr) and a non-stratified sample would lead to very few observations of deforestation events, rendering difficult a good estimation of the slope parameters for the explanatory variables. Second, only the value of the linear model intercept is affected by this balanced sampling, which is not the case for

the slope or random parameters. In our case, a biased estimate of the intercept is not an issue, as we are not interested in estimating the intensity of deforestation but the relative probability of deforestation between pixels.

Parameter inference was done in a hierarchical Bayesian framework. Non-informative priors were used for all parameters:  $\beta \sim \text{Normal}(\text{mean} = 0, \text{var} = 10^6)$  and  $V_\rho \sim 1/\text{Gamma}(\text{shape} = 0.05, \text{rate} = 0.0005)$ . We run a Markov Chain Monte Carlo (MCMC) of 7000 iterations. We discarded the first 2000 iterations (burn-in phase) and we thinned the chain each 5 iterations (to reduce autocorrelation between samples). We obtained 1000 estimates for each parameter. MCMC convergence was visually checked looking at MCMC traces and parameter posterior distributions. Function `model_binomial_iCAR()` from the `forestatrisk` Python package was used for parameter inference. This function calls an adaptive Metropolis-within-Gibbs algorithm (Rosenthal *et al.*, 2011) written in C for maximum computation speed.

## 2.3 Model comparison using percentage of deviance explained and cross-validation

We computed the deviance  $\mathcal{D}$  of the two models with the formula  $\mathcal{D} = -2\log \mathcal{L}$ ,  $\mathcal{L}$  being the likelihood of the model, i.e. the probability of observing the data given the model and estimated parameters. We compared the deviance of the two models with the deviances of both the “null” model and the “full” model. The “null” model assumes a constant probability of deforestation for all the observations and has only one parameter, the intercept of the linear relationship. At the other extreme, the “full” model has as many parameters as there are observations. We then computed the percentage of deviance explained by each model, considering that the “null” model explains 0% of the deviance and the “full” model explains 100% of the deviance.

We also performed a cross-validation to compare models using an independent validation data-set of 20,000 forest pixels in 2000. Again, the sample was stratified between 10,000 deforested pixels in 2000-2010 and 10,000 non-deforested pixels. We used the fitted models to predict the deforestation probability of all pixels in the validation data-set. To transform the deforestation probabilities into binary values, we identified the probability threshold respecting the percentage of deforested pixels (eg. the mode of the probabilities for a deforestation rate of 50%). Using model predictions and observations, we computed several accuracy indices: the Area Under the ROC Curve (AUC), the Figure of Merit (FOM), the Overall Accuracy (OA), the Expected Accuracy (EA), the Kappa of Cohen (K), the Specificity (Spe), the Sensitivity (Sen), and the True Skill Statistics (TSS). A detailed description of these indices can be found in Pontius *et al.* (2008) (for the FOM) and Liu *et al.* (2011) (for all the other indices). Formulas used to compute the indices are presented in Appendix 1.

Because the value of these indices depends on the deforestation rate (Pontius *et al.*, 2008), we computed the accuracy indices for various percentage of deforested pixels: 1, 5, 10, 25 and 50%. To do so, we selected at random subsamples of the deforested pixels in our validation data-set.

## 2.4 Forecasting deforestation

We used the fitted models to predict the deforestation probability of all the forest pixels for the year 2000. In 2000, Madagascar was covered by 9.9 Mha of natural forest corresponding to more than 109 M pixels at 30 m resolution. Predictions were computed using functions `predict_raster*`() from the `forestatrisk` Python package which make computation fast and efficient (with low memory usage) by treating raster data by blocks.

For the “icar” model, before computing the predictions of the deforestation probability, the spatial random effects at 10 km were interpolated at 1 km using a bicubic interpolation method. This was done in order to obtain spatial random effects at a resolution closer to the original 30 m resolution of the forest raster, and to smooth the deforestation probability spatially.

## 2.5 Forecasts skill scores

### 3 Results



## 4 Discussion

## 5 Acknowledgements

## 6 Tables

Table 1: **Set of explicative variables used to model the spatial probability of deforestation.** A total of height variables were tested. They described topography, forest accessibility, forest landscape, land tenure and deforestation history.

Product	Source	Variable derived	Unit	Resolution (m)
Forest maps (1990-2000- 2010)	Vieilledent et al. 2018	distance to forest edge	m	30
		distance to past deforestation	m	30
Digital Elevation Model	SRTM v4.1 CSI-CGIAR	altitude	m	90
		slope	degree	90
Highways Places	OSM-Geofabrik	distance to roads	m	150
		distance to towns	m	150
Waterways		distance to river	m	150
Protected areas	Rebioma	presence of protected area	—	30

## 7 Figures

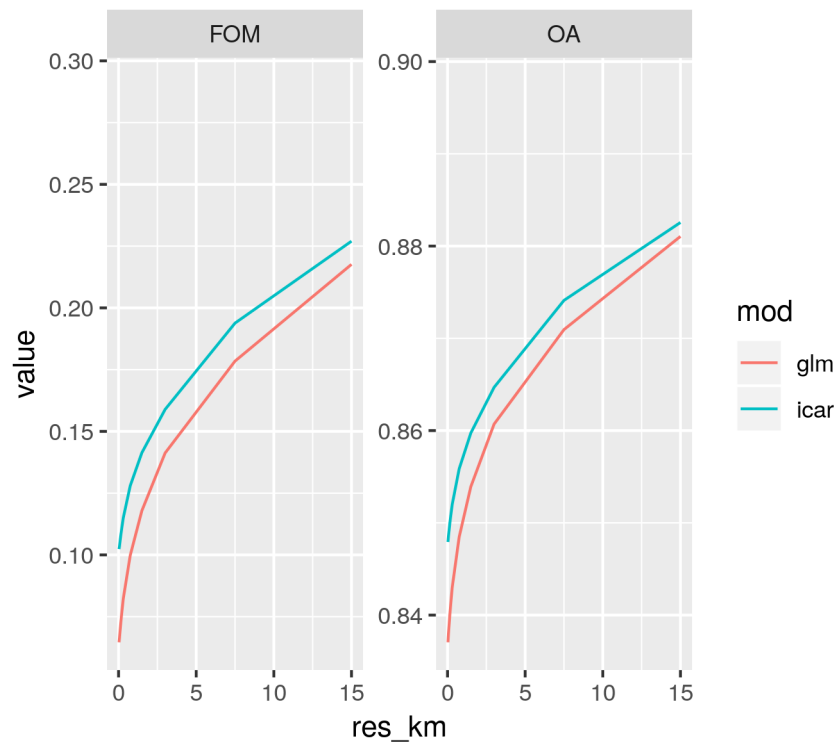


Figure 1: **Models' accuracy indices at multiple resolutions** We compared the projected deforestation maps of the two models (glm and icar) to the observed deforestation map on the period 2010-2017 (Vieilledent *et al.*, 2018a,b). We computed the Figure Of Merit (FOM) and the Overall Accuracy (OA) of the projections at multiple resolutions (from 30 m to 15 km).

## 8 Appendices

### 8.1 Appendix 1: Mathematical formulas for accuracy indices

Table 2: **Confusion matrix used to compute accuracy indices.** A confusion matrix can be computed to compare model predictions with observations.

		Observations		Total
		0 (non-deforested)	1 (deforested)	
Predictions	0	$n_{00}$	$n_{01}$	$n_{0+}$
	1	$n_{10}$	$n_{11}$	$n_{1+}$
Total		$n_{+0}$	$n_{+1}$	$n$

Table 3: **Formulas used to compute accuracy indices..** Several accuracy indices can be computed from the confusion matrix to estimate and compare models' predictive skill. We followed the definitions of Pontius *et al.* (2008) for the FOM and Liu *et al.* (2011) for the other indices. Note that the AUC relies on the predicted probabilities for observations 0 (non-deforested) and 1 (deforested), not on the confusion matrix per se.

Index	Formula
Overall Accuracy	$OA = (n_{11} + n_{00})/n$
Expected Accuracy	$EA = (n_{1+}n_{+1} + n_{0+}n_{+0})/n^2$
Figure Of Merit	$FOM = n_{11}/(n_{11} + n_{10} + n_{01})$
Sensitivity	$Sen = n_{11}/(n_{11} + n_{01})$
Specificity	$Spe = n_{00}/(n_{00} + n_{10})$
True Skill Statistics	$TSS = Sen + Spe - 1$
Cohen's Kappa	$K = (OA - EA)/(1 - EA)$
Area Under ROC Curve	$AUC = 1/(n_{+1}n_{+0}) \sum_{i=1}^{n_{+0}} \sum_{j=1}^{n_{+1}} \phi(\delta_i, \theta_j)$ where $\phi(\delta_i, \theta_j)$ equals 1 if $\theta_j > \delta_i$ , 1/2 if $\theta_j = \delta_i$ , and 0 otherwise $\delta_i$ and $\theta_j$ are the predicted probabilities for $Y_i = 0$ and $Y_j = 1$

## References

- Banerjee, S., Carlin, B.P. & Gelfand, A.E. (2014) *Hierarchical Modeling and Analysis for Spatial Data, Second Edition*. Chapman and Hall/CRC.  
URL <https://doi.org/10.1201%2Fb17115>
- Besag, J., York, J. & Mollié, A. (1991) Bayesian image restoration, with two applications in spatial statistics. *Annals of the Institute of Statistical Mathematics*, **43**, 1–20.
- Dezécache, C., Salles, J.M., Vieilledent, G. & Hérault, B. (2017) Moving forward socio-economically focused models of deforestation. *Global Change Biology*, **23**, 3484–3500.  
URL <https://doi.org/10.1111/gcb.13611>
- Liu, C., White, M. & Newell, G. (2011) Measuring and comparing the accuracy of species distribution models with presence-absence data. *Ecography*, **34**, 232–243. ISSN 1600-0587.  
URL <http://dx.doi.org/10.1111/j.1600-0587.2010.06354.x>
- Pontius, R., Boersma, W., Castella, J.C., Clarke, K., de Nijs, T., Dietzel, C., Duan, Z., Fotsing, E., Goldstein, N., Kok, K., Koomen, E., Lippitt, C., McConnell, W., Mohd Sood, A., Pijanowski, B., Pithadia, S., Sweeney, S., Trung, T., Veldkamp, A. & Verburg, P. (2008) Comparing the input, output, and validation maps for several models of land change. *The Annals of Regional Science*, **42**, 11–37. ISSN 0570-1864.  
URL <http://dx.doi.org/10.1007/s00168-007-0138-2>
- Rosenthal, J.S. *et al.* (2011) Optimal proposal distributions and adaptive MCMC. *Handbook of Markov Chain Monte Carlo*, **4**.
- Vieilledent, G., Grinand, C., Rakotomalala, F.A., Ranaivosoa, R., Rakotoarijaona, J.R., Allnutt, T.F. & Achard, F. (2018a) Output data from: Combining global tree cover loss data with historical national forest-cover maps to look at six decades of deforestation and forest fragmentation in Madagascar.  
URL <http://dx.doi.org/10.18167/DVN1/AUBRRC>
- Vieilledent, G., Grinand, C., Rakotomalala, F.A., Ranaivosoa, R., Rakotoarijaona, J.R., Allnutt, T.F. & Achard, F. (2018b) Combining global tree cover loss data with historical national forest cover maps to look at six decades of deforestation and forest fragmentation in Madagascar. *Biological Conservation*, **222**, 189–197.  
URL <https://doi.org/10.1016/j.biocon.2018.04.008>