

Dear Editor,
Please find herewith our revised ms "Environmental control of natural gap size distribution in tropical forests" for re-submission to Biogeosciences. We first would like to thank both reviewers for the high quality and seriousness of their review. We acknowledge them for that and we recognize that their works have contributed to improve the quality of our manuscript. Responses to their comments are in red in the following as well as are changes in the main ms.
We remain at your disposal for any further enquiries.
Youven Goulamoussène,

Reviewer 1

Generally, this study provides new and interesting insights in key environmental drivers of gap formation at landscape scale. While this study certainly deserves to be published, I have pointed a few issues that should be addressed before publication. The authors have developed an innovative analytic method to define gaps, but the entire analysis relies upon an a priori threshold equal to the 0.001th percentile of the estimated "natural variation of canopy height". While this choice may be well grounded, the rationale beyond it remains unexposed. How sensitive is the definition of gap and all subsequent results to this threshold? What if the authors had chosen the 0.01th percentile? Some kind of sensitivity analysis would make their choice more reliable.

We easily admit that our initial choice of the 0.001th percentile for the height threshold may have been seen as arbitrary. We choose this threshold value to keep the maximum of information from the first distribution (gap height) while minimizing biases due to including too much information from the 2nd one (canopy height). Because the whole process needs several hours to be performed, from the GIS works to the model inference, with thousands of gaps, we cannot run a well-performed sensitivity analysis. However, we provide in supplementary informations, the parameter values for 2 additional thresholds (0.0001th and 0.01th). The posterior values of almost all variables are quite similar and thus do not change interpretation: Slope, TRI are always positive whatever the threshold TOPEX, HAND are always negative whatever the threshold DA and HAlt always include zero in the credibility interval

More importantly, while a landscape scale approach seems meaningful to infer gap size distribution, this study highlights the importance of environmental factors on both gap frequency and size. Thus, I wonder if a fixed definition of gap remains meaningful, or if that definition should not adapt to the different forest types and/or main topographical features found at large scale. Doing so would point towards a more "ecological" definition of gaps, instead of a pure remote-sensing approach, and ultimately raises the question of the aim of detecting gaps. For instance, does a 100m² gap in waterlogged areas dominated by *Euterpe oleracea* has the same ecological meaning than on hilltops? Certainly not in term of number of trees killed, biomass loss and forest turnover. Depending on the variable of interest (e.g. carbon emission), a fit-them-all definition is questionable. This issue, if not formally addressed, should at least be discussed. The following recent publications may provide additional information (Chambers et al., 2013; Lobo and Dalling, 2014; Schliemann and Bockheim, 2011).

The choice of the values of height and threshold may be adapted to different forest types and topographic characteristics. In our case, the choice was fully data-driven using the DCM and DEM and no ecological knowledge. Within our framework it is likely that in waterlogged areas, areas covered with mature trees that do not exceed the height thresholds may appear in our analysis as forest gaps. In order to clarify this question, an approach using time-series would allow to identify these 'false' gaps that never get filled and thus are not part of the forest endogenous dynamics. These are not gaps in the ecological meaning.

Finally, the manuscript requires additional efforts in editing (loads of typo citations errors, unclear headers and acronyms) and reviewing recent literature (lots of relevant publications is missing, comparing lambda with other studies). A proof-reading by a native English would also greatly help. This has been done. The manuscript has been edited by a professional science editing service. We do

believe that directly comparing lambda values between studies is difficult because it may depend on the assumed (or inferred) height and size thresholds.

Specific comments

- l. 99 : For clarity, please define explicitly all the acronyms used, i.e. topographic exposure (TOPEX).
See line 99
- l.119 : Sub-header should be : “Height above the nearest drainage” to be consistent with previous sub-headers
See line 120 . Section Methods 2.1.2.
- l.193 : What is the resolution of the TOPEX variable? Do you have several indices by 5m2? Please clarify how you can get 2 values (min max), or did you standardize TOPEX as : $\text{abs}((\text{TOPEX} - \min(\text{TOPEX})) / (\max(\text{TOPEX}) - \min(\text{TOPEX})))$.
The native pixel resolution is 5 m \times 5m. Original values for TOPEX are hardly interpretable because they are counter-intuitive. In order to simplify the interpretation, we thus modified the TOPEX values to get the highest for the highest exposure. See the variable transformation section.
- l. 212 : I suggest to change the header here, as Kuo-Mallik refers to a method, but you used it to select the variables. “Variables selection” looks more appropriate. There is also an issue in the way the reference is quoted.
Line 212 : Agreed
- l. 216 : there is an missing (or extra) parenthesis in your expression
Line 217 : done
- l. 226 : “Given this height, we retained the surface $x_{\min} = 104 \text{ m}^2$ ”. What is the link between the height threshold and the minimal gap area, here? I thought both minimal height and gap size were defined separately.
Yes they are. Firstly, we define a height threshold from which we observe a gap. Secondly we used this height threshold in order to determine the minimum gap size area using the Kolmogorov-Smirnov (KS) distance criterion
- l.254 : "Environmental covariates with posterior KM values close to 1 , NAMELY Slope, TOPEX, and HAND ...
Line 254 : Done. "Environmental covariates with posterior KM values close to 1, namely Slope, Topex, and HAND"
- l.260 : “Defining the height threshold at which forest gaps may be delineated is a major difficulty faced by foresters. Many times, canopy gaps have been defined in the field, adopting Brokaw’s definition” Is it only the minimal height that is at stake here, or also the minimal area? Many studies define gaps regarding to their size (e.g. Denslow et al., 1998 ; Hérault et al., 2010 ; Lima, 2004). This sentence sounds odd, I suggest rephrasing as follow : “Delineating forest gaps is a persistent challenge for foresters and ecologists, among whom Brokaw’s gap definition (1982) has remained very popular/extensively used.”
We thank you for your suggestion. The changes were made at line 263
- l. 265-269 : There are several studies that do not use this 2m-threshold definition of gaps, but 10m (e.g. Hubbell et al., 1999 ; Meer and Bongers, 1996 ; Welden et al., 1991). While the authors are extensively referring to the seminal paper of Brokaw, there are way more references defining gaps in complex tropical forests that are lacking here.
We have updated the bibliography
- l.300 (onwards) : This paragraph is very confusing. Where does the 75% comes from? What is the remain 25% then? In sloppy areas, does it make a big difference if a tree falls due to

breakage, or being uprooted? I don't think so, and tree size seems to be a more important factor in the cascading effect than mode of death. Yet, the turnover may be more rapid on slopes than bottomland, resulting in fewer large trees (and tree fall gaps).

Rewritten

— 1.338 : but WE found

Done, line 345

1 Supplementary information

TABLE 1 – List of environmental variables, abbreviations, units, and values of the posteriors in univariate models for a height threshold equal to the 0.0001th percentile of the height distribution of the canopy.

Parameter	Abbreviation	Unit	Posterior value	Confidence interval (CI 95%)
Slope	Slope	°	0.119	[0.0416 ; 0.208]
Terrain Ruggedness Index	TRI	-	0.119	[0.083 ; 0.157]
TOPographic EXposure	TOPEX	-	-0.128	[-0.188 ; 0.00202]
Drained Area	DA	m ²	0.0843	[-0.0574 ; 0.179]
The Hydraulic Altitude	HAIt	m	-0.0135	[-0.04 ; 0.042]
HAND	HAND	-	-0.0615	[-0.152 ; 0.0162]

TABLE 2 – List of environmental variables, abbreviations, units, and values of the posteriors in univariate models for a height threshold equal to the 0.001th percentile of the height distribution of the canopy.

Parameter	Abbreviation	Unit	Posterior value	Confidence interval (CI 95%)
Slope	Slope	°	0.0735	[-0.02 ; 0.15]
Terrain Ruggedness Index	TRI	-	0.0718	[0.04 ; 0.10]
TOPographic EXposure	TOPEX	-	-0.082	[-0.12 ; -0.05]
Drained Area	DA	m ²	-0.0176	[-0.09 ; 0.05]
The Hydraulic Altitude	HAIt	m	-0.0177	[-0.05 ; 0.02]
HAND	HAND	-	-0.003	[-0.08 ; 0.09]

TABLE 3 – List of environmental variables, abbreviations, units, and values of the posteriors in univariate models for a height threshold equal to the 0.01th percentile of the height distribution of the canopy.

Parameter	Abbreviation	Unit	Posterior value	Confidence interval (CI 95%)
Slope	Slope	°	0.0975	[-0.02 ; 0.17]
Terrain Ruggedness Index	TRI	-	0.089	[0.05 ; 0.12]
TOPographic EXposure	TOPEX	-	-0.012	[-0.03 ; -0.32]
Drained Area	DA	m ²	-0.004	[-0.08 ; 0.05]
The Hydraulic Altitude	HAIt	m	0.063	[-0.04 ; 0.08]
HAND	HAND	-	-0.01	[-0.09 ; 0.06]

Environmental control of natural gap size distribution in tropical forests

Youven Goulamoussène^{1,2}, Caroline Bedeau³, Laurent Descroix³,
Laurent Linguet², and Bruno Hérault¹

¹Cirad, UMR EcoFoG (AgroParisTech, CNRS, Inra, Univ Antilles, Univ Guyane), Kourou, French Guiana

²Université de Guyane – UMR Espace-Dev, BP 792, 97337 Cayenne, France

³Office National des Forêts (ONF), département Département RD, Cayenne, French Guiana

Correspondence to: Youven Goulamoussène (youven.goulamoussene@ecofog.gf)

Abstract. Natural disturbances are the dominant form of forest regeneration and dynamics in unmanaged tropical forests. Monitoring the size distribution of treefall gaps is important to better understand and predict the carbon budget in response to land use and other global changes. In this study, we model the size frequency distribution of natural canopy gaps with a discrete power law distribution. We use a Bayesian framework to introduce and test, using Monte Carlo Markov Chain and Kuo-Mallick algorithms, the effect of local physical environment on gap size distribution. We apply our methodological framework to an original Light Detecting and Ranging dataset in which natural forest gaps were delineated over 30000 ha of unmanaged forest. **We highlight strong links between gap size distribution and environment, primarily hydrological conditions and topography, with large gaps being more frequent in floodplains and on wind-exposed areas.** In the future, we plan to scale up testing our methodology with satellite data. Additionally, although gap size distribution variation is clearly under environmental control, gap process variation over time should be tested against climate variability.

1 Introduction

Natural disturbances caused by forest gaps play an important role in tropical rainforest dynamics. Canopy gaps caused by the death of one or more trees are the dominant form of forest regeneration because the creation of canopy openings continuously reshapes forest structure as gaps are filled with younger trees (Whitmore, 1989). The first, and perhaps most important, effect of gap occurrence is an immediate increase in light intensity (Hubbell et al., 1999a), allowing sunlight to penetrate the understory. This phenomenon has been widely studied because the opening of gaps contributes

to the establishment and growth of light-demanding trees (Denslow et al., 1998), thus contributing to the maintenance of biodiversity. Another effect of canopy gaps is the local modification of the forest nutrient balance (Rüger et al., 2009). **When canopy gaps are created, large amounts of dead leaves and wood will be decomposed and mineralised so that the availability of soil nutrients for neighboring trees will increase** (Brokaw, 1985). These nutrient patches are also linked to small-scale spatial variations in forest carbon balance, as shown by Feeley et al. (2007). The relationship between gap formation and the population dynamics of trees or lianas is also quite well understood, with increased liana basal area (Schnitzer et al., 2014) allowing low-wood-density pioneer species to recruit exclusively in newly formed gaps (Molino and Sabatier, 2001).

Many studies have investigated the effect of treefall gaps on biodiversity, particularly animal movement and species composition (Bicknell et al., 2014; Puerta-Piñero et al., 2013), carbon cycles, and forest dynamics. Some authors use field data to study natural gap dynamics, usually at plot scale (Hubbell et al., 1999b). As these studies are quite limited in spatial extent (< 50 ha) and because gap formation is largely unpredictable (Hubbell et al., 1999a; Lloyd et al., 2009), optical satellite imagery has been widely promoted and proven adequate for monitoring forest gaps over space and time (Frolking et al., 2009). At high resolution (< 10 m), IKONOS satellite images may be well suited for evaluating gap dynamics (Espírito-Santo et al., 2014). In French Guiana, the SPOT-4 satellite (20 m spatial resolution) has successfully detected canopy gaps (Colson et al., 2006) using a combination of several spectral bands, such as near and short-wave infrared. However, topographical variation, gap shape, and shade may influence and bias gap detection with optical products. Moreover, persistent cloud cover, common in many tropical forest basins, limits their utility.

Airborne Light Detecting and Ranging (LiDAR) platforms therefore offer a solution to this problem. Recent developments in LiDAR have significantly advanced our ability to derive accurate measurements of canopy forest structure, to detect gaps, and to assess the effect of spatial and temporal variation in carbon balance (Asner and Mascaro, 2014). Kellner and Asner (2009) used remote LiDAR sensing to quantify canopy height and gap size distributions in five tropical rain forest landscapes in Costa Rica and Hawaii. They showed that canopy gaps can be observed with the help of LiDAR-derived digital canopy models (DCMs) and that gap size frequency distribution (GSFD) can be fit with a power law distribution, suggesting a surprising similarity in canopy gap size frequency distributions on diverse soil types associated with diverse geologic substrate ages. Asner et al. (2013) also used LiDAR data to analyze whether gap size frequency distribution is modified by topographic and geologic characteristics and again showed that canopy gap size distribution is largely invariant between forests on erosional terra firme and depositional floodplain substrates in the Peruvian Amazon basin. Finally, using airborne LiDAR, Lobo and Dalling (2014) have recently explored the effect of forest age, topography, and soil type on canopy disturbance patterns across central Panama. For the first time, they highlighted significant effects of slope and of forest age, with a higher frequency of large gaps associated with old-growth forests and gentle slopes.

In this study, we use a DCM derived from airborne LiDAR across a 30000 ha tropical forest landscape in the Régina forest in French Guiana. This approach provides high-resolution maps of canopy gaps and helps us to understand the environmental determinism of gap occurrence in tropical forests. Our specific aims were therefore:

- to define canopy gaps from canopy height data using a probabilistic approach
- to model gap size distribution by inferring a likelihood-explicit discrete power law distribution in a Bayesian framework
- to introduce the environment into the scaling parameter of the power law distribution and test its predictive ability

2 Materials and Methods

The study site is located in the Régina forest (4°N, 52°W), where the most common soils are ferallitic. The site is located on slightly contrasting plateau-type reliefs that are rarely higher than 150 m on average. The forest is typical of French Guianese rainforests. Dominant plant families in the Régina forest include *Burseraceae*, *Mimosoideae*, and *Caesalpinoideae*. The site receives 3,806 mm of precipitation per year, with a long dry season from mid-August to mid-November, and a short dry season in March (Wagner et al., 2011).

2.1 Data source

2.1.1 LiDAR data

LiDAR data were acquired by aircraft in 2013 over 30,000 ha of forest by a private contractor, Altoa (<http://www.altoa.fr/>), using a Riegl LMS-Q560 laser. This system was composed of a scanning laser altimeter with a rotating mirror; a GPS receiver (coupled to a second GPS receiver on the ground); and an inertial measurement unit to record the pitch, roll, and heading of the aircraft. The laser wavelength was near-infrared (from about 800 nm to 2500 nm). Flights were conducted at 500 m above ground level with a ground speed of 180 km.h⁻¹, and each flight derived two acquisitions. The LiDAR was operated with a scanning angle of 60° and a 200 kHz pulse repetition frequency. The laser recorded the last reflected pulse with a precision better than 0.10 m, with a density of 5 pulses.m⁻².

The DCM was derived from the raw scatter plot consisting of the pooled dataset from the two acquisitions. Raw data points were first processed to extract ground points using the TerraScan (TerraSolid, Helsinki) ground routine, which classifies ground points by iteratively building a triangulated surface model. Ground points typically made up less than 1% of the total number of the return pulses. The DCM has a resolution of 1 m. In order to remove areas close to natural rivers, a 20 m buffer was first applied to all shorelines. Then a 25 m buffer was applied to anthropogenic tracks.

2.1.2 Environmental data

We use six environmental variables to synthesize the observed environmental gradients. All variables were computed from a LiDAR digital terrain model (DTM) with 5 m² cells.

Slope

- 95 The slope was derived from the LiDAR DTM. Slope was computed at a grid cell as the maximum rate of change in elevation from that cell to its 8 neighboring cells over the distance between them.

Topographic exposure

- We use the **TOPographic EXposure** (TOPEX) index to measure topographic exposure to wind (Chapman, 2000). TOPEX is a variable that represents the degree of shelter assigned to a location. It was
100 derived from quantitative assessment of horizon inclination. The values of this index are closely correlated with wind-shape index (Mikita and Klimánek, 2012). Exposure is calculated based on the height and distance of the surrounding horizon, which are combined to obtain the inflection angle. We use this angle to quantify topographic exposure (**pixel resolution 5 m × 5m**). When a large topographic feature, like a mountain, is far off in the distance the inflection angle is low. When the
105 same mountain is closer, the inflection angle is higher. Therefore, a higher inflection angle is equal to lower exposure or higher sheltering (Mikita and Klimánek, 2012).

Drained area

- Drained area (DA) measures the surface of the hydraulic basin that flows through a cell. A low value indicates that a cell is located at the border between two basins, whereas high values indicate cells
110 located downstream.

Hydraulic altitude

The hydraulic altitude (HA) of each cell, its altitude above the closest stream of its hydraulic basin, was computed from the 3rd order hydraulic system. Low values, including 0, indicate that the forest plot is potentially temporarily flooded, whereas high values indicate that it is located on a hilltop.

115 Terrain ruggedness index

The terrain ruggedness index (TRI) captures the difference between flat and mountainous landscapes. TRI was calculated using SAGA GIS SAGA (2013) as the sum of the altitude change between a pixel and its eight neighboring pixels (Riley, 1999).

The height above the nearest drainage

120 The height above the nearest drainage (HAND) model normalizes topography with respect to drainage network by applying two procedures to the DTM. The initial basis for the HAND model came from the definition of a drainage channel: perennial streamflow occurs at the surface, where the soil substrate is permanently saturated. It follows that the terrain at and around a flowing stream must be permanently saturated, independently of the height above sea level at which the channel occurs.

125 Streamflow indicates the localized occurrence of homogeneously saturated soils across the landscape. The second basis for the HAND model came from the distinctive physical features of water circulation. Land flows proceed from the land to the sea in two phases: in restrained flows at the hillslope surface and subsurface, and in freer flows (or discharge) along defined natural channels. (Nobre et al., 2011)

130 2.2 Forest gap definition

2.2.1 Height threshold

To identify discrete canopy gaps, we had to choose a gap threshold height. Some authors define this threshold at 2 m (Brokaw, 1982). Runkle (1982) defines a gap as the ground area under a canopy opening that extends to the base of the surrounding canopy trees, these usually being considered to

135 be taller than 10 m, with a trunk diameter at breast height (DBH) > 20 cm. However, in practice, defining gap boundaries is a tricky issue, even in the field. Here, we develop a probabilistic method for detecting canopy gaps from LiDAR data. We used the DCM to model canopy height distribution considering a mixture distribution of two ecological states: the natural variation of canopy height in mature forests, modeled as a normal distribution, and the presence of forest gaps, which lead to a

140 new normal distribution with lower values. We consider that the threshold between the two states is equal to the 0.001th percentile of the height distribution of the canopy (our results appeared robust to the threshold value, see supplementary information). We then define canopy gaps as contiguous pixels at which the vegetation height is less than or equal to the height threshold. Contiguous pixels are defined as pixels that have any contact contact by edges or by vertices.

145 2.2.2 Minimum gap size

In our study, we define the minimum area of a gap as x_{min} . We model the gap size frequency distribution with a power law distribution. We use the Pareto distribution in a discrete power law probability density function (Virkar and Clauset, 2014). These distributions have a negative slope and their size frequencies are plotted on logarithmic axes, allowing us to observe the scaling parameter λ .

150 A value close to 1 means there are a large number of large gaps. In other word, in forests dominated by small canopy openings, values of λ are larger, whereas smaller values of λ increase the frequency

of large gap events (Fisher et al., 2008). In a discrete power law with parameter λ , the probability for gap size x is given by:

$$p(x) = \frac{x^{-\lambda}}{\zeta(x_{min}, \lambda)}, \quad (1)$$

155 where x_{min} is the lower truncation point and λ is the scaling parameter.

Most of the analysis was performed under R (Team et al., 2013) and making use of `powerLaw` (Clauset et al., 2009) and `VGAM` (Yee et al., 2010) packages.

We use a Kolmogorov-Smirnov (KS) distance criterion order to determine the error between the observed distribution and the Pareto distribution. KS is defined as the maximum distance between
160 the cumulative distribution functions (CDFs) of the data and the fitted function (Virkar and Clauset, 2014). We retain, for the remainder of this study, a minimum gap size area $x_{min} = 104 \text{ m}^2$, which minimized the KS distance in our dataset.

2.3 Modeling gap size distribution

Having set the height threshold and minimum gap size, the GSFD is modeled with a discrete Pareto
165 distribution frequency.

2.3.1 Model inference

We use a Bayesian framework to estimate model parameters. Here, the value of a parameter is estimated by its posterior distribution, which by definition, is proportional to the product of the likelihood of the model and the parameter prior distribution. The prior distribution is based on prior
170 knowledge of the possible values of a parameter. The posterior densities of the different parameters were estimated using a Monte Carlo Markov Chain algorithm (MCMC).

2.3.2 Metropolis-Hastings algorithm

As the model contains many parameters, we built a Metropolis-Hastings (MH) algorithm in which all parameters are updated together. Details on the algorithm are given below:

- 175 • $Y = y_1, y_2, \dots, y_n$ is the gap size vector
- $X = x_{g1}, x_{g2}, \dots, x_{gi}$ is the vector of covariates (environmental variables) for gap g
- $\theta = \theta_1, \theta_2, \dots, \theta_i$ is the model parameter vector

The first values of the parameter vector are initialized as $t = 1, \theta^t \sim \pi_{\theta}^0$.

180 For each step t , a new parameter value is sampled from the proposition distribution and a new vector of theta candidates is generated.

$$\theta^{cand} \sim \pi^{prop} \quad (2)$$

Acceptance or rejection of the new candidate θ^{cand} is determined by computing the likelihood ratio of the two discrete Pareto distributions:

$$185 \quad \rho(\theta^t, \theta^{cand}) = \underbrace{\frac{\mathcal{L}(Y|X, \theta^{cand})}{\mathcal{L}(Y|X, \theta^t)}}_{\text{likelihood}} \underbrace{\frac{\pi_{\theta}^0(\theta^{cand})}{\pi_{\theta}^0(\theta^t)}}_{\text{prior}} \underbrace{\frac{\pi^{prop}(\theta^t)}{\pi^{prop}(\theta^{cand})}}_{\text{proposal}} \quad (3)$$

The candidate θ^{cand} is accepted or rejected as follows:

$$u \sim \mathcal{U}_{[0,1]}, \theta^{cand} \begin{cases} \theta^{t+1} & \text{if } u < \rho(\theta^t, \theta^{cand}) \\ \theta^t & \text{if } u > \rho(\theta^t, \theta^{cand}) \end{cases} \quad (4)$$

The algorithm is run for 1000 iterations. We use the median of the posterior densities to estimate parameter values, and the distribution of the posterior densities to estimate parameter credibility
190 intervals.

2.3.3 Univariate environmental effects

Variable transformation

To improve model inference, parameter significance and interpretation, we first transformed some environmental variables:

$$195 \quad Slope = \text{sqrt}(slope) \quad (5)$$

$$HAlt = \log(HA + 1) \quad (6)$$

$$TopeX = |\max(TOPEX) - (TOPEX)| \quad (7)$$

The environmental variables are then centered and scaled with R function : "scale".

We first consider each environmental covariate independently. These covariates are included one-
200 by-one in the model to constrain the exponent λ . We use the exponential function to constrain λ , because the Riemann's zeta function **only admits $\lambda > 1$** .

$$\lambda_i = 1 + \exp(\theta_{0i} + \theta_i \times X) \quad (8)$$

where λ_{ig} is the λ value dependent on the value of environmental variable i in gap g , θ_{0i} is the intercept, and θ_i quantifies the effect of covariates var on the gap size distribution var_{ig} .

205 2.3.4 Multivariate model

Principal component analysis

We first investigated the collinearity of environmental data through principal component analysis (PCA) on the normalized environmental dataset.

Model

210 To build the final model, we used the results of the univariate model (Table 1) and the PCA (Figure 3) and set:

$$\lambda = 1 + \exp(\theta_0 + \theta_1 \times Slope + \theta_2 \times Topea + \theta_3 \times HAlt + \theta_4 \times HAND) \quad (9)$$

Variables selection

215 To select the significant covariates and build the final model, we used the method proposed by **Kuo and Mallick (1998) (KM)**. This method consists of associating an indicator with each variable var_i and parameter θ_i . This indicator can take two values: 1 or 0. If it is set to 1, the variable is included in the model, but if the value is set to 0, it is not. We used the MH and KM algorithms to estimate the indicators I and infer their *a posteriori* distribution in addition to θ .

220 We start the KM algorithm with $t = 1$, $\theta^t \sim \pi_\theta^0$, $I_j^t \sim Ber(0.5)$ for $j = 1, \dots, i$. For each covariate j (selected in random order), we use the MH algorithm to update θ_j . To update I_j , we compute the ratio ρ (eq10) and generate I_j^{t+1} from a Bernoulli distribution $Bern(\rho)$:

$$\rho = \frac{1}{1 + \frac{\mathcal{L}(Y|\tilde{X}, \theta^t, I_j=0, I_{-j}^t)}{\mathcal{L}(Y|\tilde{X}, \theta^t, I_j=1, I_{-j}^t)}} \quad (10)$$

Model inference and data analysis were conducted with R software (R-Core-Team 2012). All
225 maps and geographical information were computed with SAGA (SAGA, 2013) and ArcGIS 10.1.

3 Results

3.1 Gap delineation

In this study, we used a forest canopy height mixture model to define the maximum height of a given pixel to be included in a forest gap. This probabilistic method produced results that fit the observed
230 canopy height distribution. We retained the 11 m threshold that corresponds to the 0.001th percentile of the canopy height distribution (Figure 1). Given this height, we retained the surface $x_{\min} = 104 \text{ m}^2$ that minimized the KS distance between predictions and observations. Here, our gap definition was therefore defined as an area $> 104 \text{ m}^2$, in which the LiDAR measured canopy height is always $\leq 11 \text{ m}$.

235 3.2 Basic statistics

We mapped 12,293 gaps with vegetation $\leq 11 \text{ m}$ in height. The mean gap size was 236 m^2 with a minimum gap size of 104 m^2 and a maximum of $29,063 \text{ m}^2$. The total gap area was about 290 ha,

or 1% of the whole surveyed area. The observed gap size distribution was modeled with a Pareto distribution (Figure 2), leading to a scaling parameter $\lambda_{x_{min}}$ of 2.6.

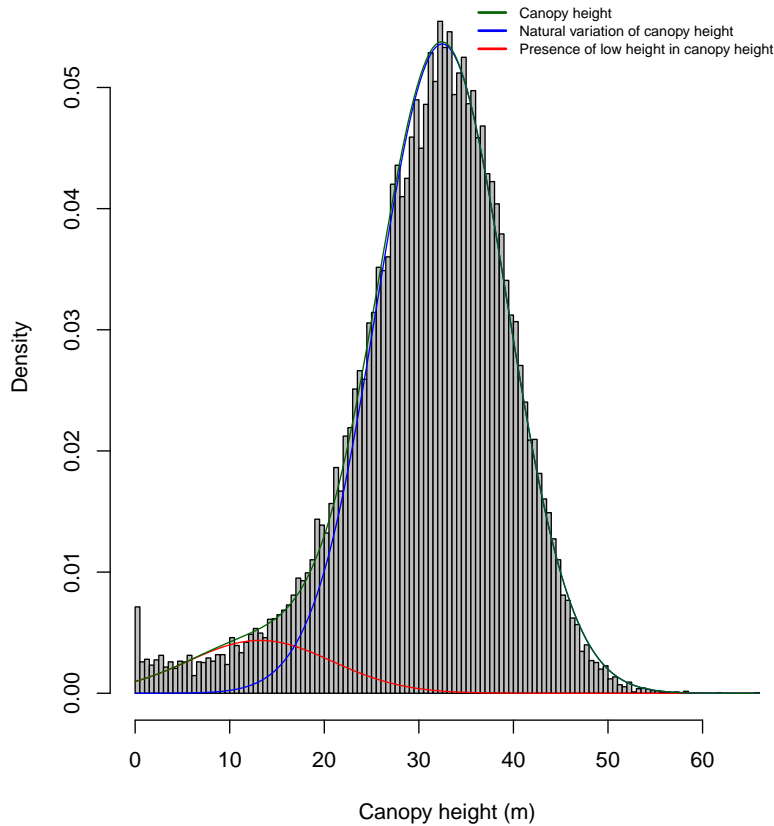


Figure 1. Canopy height distribution. Canopy height considered as a mixture distribution of two ecological features. The first (blue curve) is the natural variation in canopy height, modeled as a normal distribution. The second (red curve) is linked to the presence of low heights in the total canopy height distribution, likely to be due to a forest gap. We set the gap threshold to the 0.001th percentile of the blue curve density, *i.e.*, 11 m.

240 3.3 Univariate models

All variables had an effect on gap size distribution (Table 1). The scaling coefficient λ is related to the ratio of small gaps to large gaps, with values close to 1 indicating a higher frequency of large gaps and vice versa. Parameter estimates for slope and TRI show high occurrence of small gaps for large values of the two variables. Contrarily, the effect of DA, HAND, HAlt, and Topex on λ are
245 clearly negative, meaning that the frequency of large gaps increases with large values.

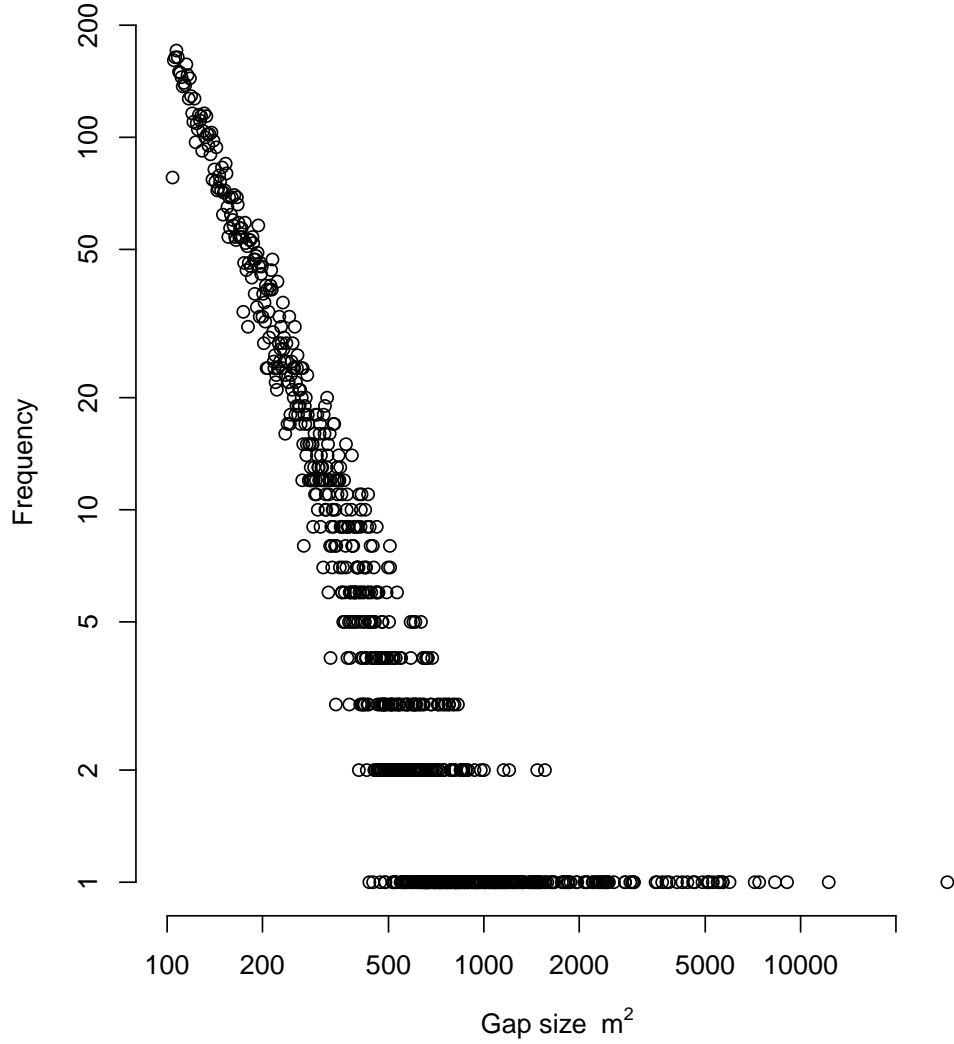


Figure 2. The observed gap size frequency distributions modeled as a power law function with $\lambda = 2.6$.

3.4 The multivariate model

To define the final multivariate predictive model, we used the significant results of the univariate models together with the output of the PCA, in order to avoid multicollinearity.

3.4.1 Variable selection

250 The first three PCA axes explained more than 80% of the data variance. The first axis, which accounted for 36.45% of the variance, was positively correlated with relative HAlt and negatively correlated to HAND and DA, and thus clearly highlighted the local altitudinal gradient. The second

Table 1. List of environmental variables, abbreviations, units, and values of the posteriors in univariate models.

Parameter	Abbreviation	Unit	Posterior value	Confidence interval (CI 95%)
Slope	Slope	°	0.0735	[-0.02 ; 0.15]
Terrain Ruggedness Index	TRI	-	0.0718	[0.04 ; 0.10]
TOPOgraphic EXposure	TOPEX	-	-0.082	[-0.12 ; -0.05]
Drained Area	DA	m ²	-0.0176	[-0.09 ; 0.05]
The Hydraulic Altitude	HA	m	-0.0177	[-0.05 ; 0.02]
HAND	HAND	-	-0.003	[-0.08 ; 0.09]

axis explained an additional 28.5% of variance and was positively correlated with the TRI and Slope. The third axis explained a further 15.2% of the variance and was correlated only with Topex (Figure 3). The multivariate model was created using a Bayesian framework including four environmental variables: slope, Topex, HAND, and HAlt, the explanatory variables that had an effect on λ . Finally, the KM methodological framework was used to select the most parsimonious model.

Environmental covariates with posterior KM values close to 1, namely Slope, Topex, and HAND (eqn 9) were retained in the final model (Figure 4). Parameter estimates of the final model indicated that the greatest effects on gap size distribution were caused by Topex and HAND.

4 Discussion

4.1 Methodology

4.1.1 Gap Detection

Delineating forest gaps is a persistent challenge for foresters and ecologists, among whom Brokaw (1982) gap definition has remained extensively used, in which "a 'hole' in the forest extending through all levels down to an average height of 2 m above the ground," must be defined by an experienced observer. There are several studies that do not use this 2m-threshold definition of gaps, for instance 10 m (e.g. Hubbell et al., 1999; Meer and Bongers, 1996; Welden et al., 1991). However, in this study we have decided to use a probabilistic approach, modeling height distribution as a mixture of two normal laws. We found a height, 11 m, which is much higher than that in Brokaw's definition, but is consistent with our field experience, where woody debris, dead canopy tree boles, and residual saplings (*i.e.*, remnants that survive the gap formation event) may rise well above 2 m. For example, Hubbell et al. (1999a) showed that small stems frequently remained in gaps up to 4-5 m in height, while Lieberman et al. (1985) reported broken and damaged stems up to 10 m tall within a gap. The choice of the values of height and threshold may be adapted to different forest types and

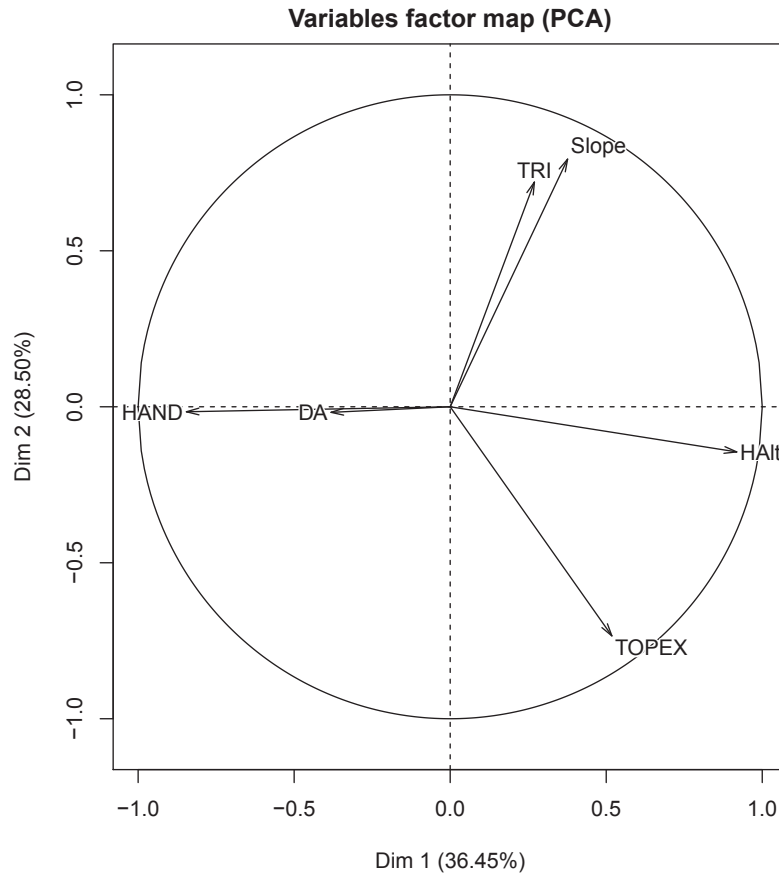


Figure 3. Results of the principal component analysis on the environmental variables

topographic characteristics. In our case, the choice was fully data-driven using the DCM and DEM and no ecological knowledge. Within our framework it is likely that in waterlogged areas, areas covered with mature trees that do not exceed the height thresholds may appear in our analysis as forest gaps. In order to clarify this question, an approach using time-series would allow to identify

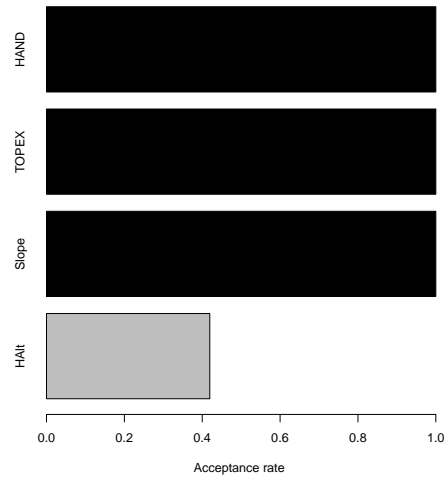


Figure 4. Results of the Kuo-Mallick algorithm for variable selection. Variables were included in the final model when their value was close to 100%: *Slope*, *Topex* and *HAND*

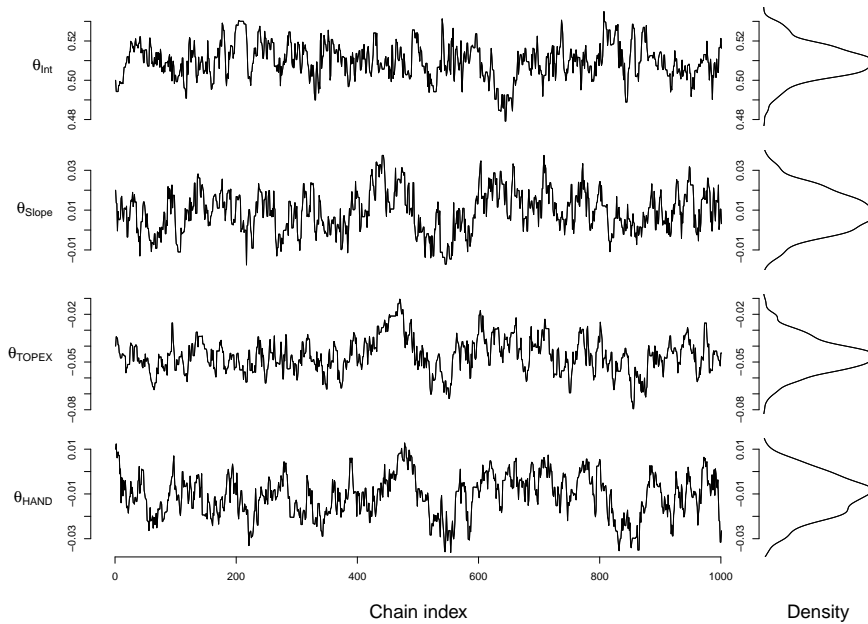


Figure 5. Posterior distribution of the environmental variables in the final multivariate model.

280 these 'false' gaps that never get filled and thus are not part of the forest endogenous dynamics.
 These are not gaps in the ecological meaning.

Defining minimum gap size is also a delicate proposition. Some authors, working with high-definition LiDAR data, have considered a minimum gap size (x_{min}) of 1 m² (Asner et al., 2013) (Kellner and Asner, 2009). This minimum gap size is unrealistic from an ecological perspective given that a hole of several square meters in the canopy may simply reflect the distance between two crowns. Brokaw recommended a range from 20 m² to 40 m² based on his field experience. We have worked with a minimum gap size of 104 m², and based this value on the minimized Kolmogorov-Smirnov distance between observed and predicted values.

We built on previous studies that show that gap size distribution follows a power law distribution. However, the underlying mechanisms that control this distribution are still unclear. The Bayesian framework we developed allowed us to detail the contributions of each environmental variable to the size of each individual gap. Because the precise environmental variables were explicitly taken into account in the model likelihood of each gap, we were able to predict gap size distribution from environmental covariates, a difficult task when the scale exponent is estimated once, at the forest level, and compared between forests. The global scale exponent that we estimated for an average environment ($\lambda = 2.6$), is consistent with some previous studies (Kellner and Asner, 2009; Kellner et al., 2011), though slightly larger than those of others Lobo and Dalling (2014) [1.97 ; 2.15] and Asner et al. (2013) [1.70 ; 2.03].

4.2 Environmental effects on gap size frequency distribution

For the first time, gap size distribution integrates environmental variables as a linear combination of the scale parameter (λ) of a discrete Pareto distribution frequency. Our results suggest that three covariates drive the gap size frequency distribution in our forest: *Slope*, *HAND*, and *TOPEX* (Figure 5).

4.2.1 Slope

Steep slopes are well-known to directly impact tropical forest canopy structure (Bianchini et al., 2010). In this study, we found similar results to Lobo and Dalling (2014) in BCI, *i.e.*, large gaps (smallest λ) are more frequent on the lowest slopes. This may seem counter-intuitive at first, as treefall may be (i) more prone to induce cascading effects when slopes are steep and (ii) more frequent in slopes where soils are shallow with lateral drainage (Gourlet-Fleury et al., 2004), impeding deep rooting of trees. However, the forest turnover is more important in bottomlands where slopes are gentle (Durrieu de Madron, 1994). Considering that large gaps may be created solely by contiguous and independent treefalls, larger gaps may then be expected in bottomlands from a pure probabilistic approach. And given the positive link between wood density and steep slopes (Ferry et al., 2010), trees may be more resistant to cascading effects than they are in bottomlands.

315 4.2.2 Water Saturation

HAND is a binary variable that takes the value 1 on water-saturated soils. Because λ decreases when *HAND* equals 1, the frequency of large gaps increases in floodplains and bottomlands. These results support the findings of (Korning and Balslev, 1994), highlighting more dynamic forests in floodplains subject to large flooding events that lead to cascading treefall events. Together with
320 (Asner et al., 2013), our results suggest that we can effectively extend these results to bottomlands, where we already know that aboveground biomass and mean wood density are 10% lower than on hilltops (Ferry et al., 2010). Given its ease of implementation on a land-surface model and its high predictive power, *HAND* covariates present great potential applicability for gap size distribution prediction.

325 4.2.3 Topographic Exposure

The effect of topographic exposure on λ is consistent with our *a priori* hypothesis that wind-exposed areas would have a greater relative frequency of large gaps. Although hurricane damage does not occur in continental equatorial regions of the Amazon (Nelson et al., 1994), here we demonstrate that tree exposure has a large impact on gap size distribution. Lobo and Dalling (2014) observed no clear
330 effect of TOPEX, and suggested that this index has a slight negative effect on gap size distribution. The results of this study are in line with the pioneering work of (Negrón-Juárez et al., 2014), which showed that wind exposure is related to higher elevations that inflate the occurrence of larger gaps. However, coastal French Guianese forests exhibit different landscapes and landforms (Guitet et al., 2013). Our study area is made of dissected plateaus characterized by simple forms resembling hills
335 (Guitet et al., 2013). It is possible that these characteristics, leading to unique combinations of landform elevations, may create complex terrain interactions that increase wind local speed and, in turn, cause large gaps. We conclude that topographic exposure is an appropriate index for predicting gap size distribution, but this must be confirmed in other landscape types.

5 Conclusions

340 To our knowledge, this is the first study where the precise environmental descriptors associated to each canopy gap were explicitly taken into account in the general model likelihood. We were able to do so because we wrote the general model likelihood as the product of all the single likelihoods (*i.e.* each gap had its own likelihood depending on the environmental covariate values). Doing so, we were able to predict gap size distribution from the fine environmental covariates, an impractical
345 task when the scale exponent is estimated once at the forest level (*i.e.* mixing all the found gaps together) and compared between forests *a posteriori*. We also put forward an innovative method to define a height threshold and minimum gap size using two probabilistic approaches. The modeled distribution of canopy height as mixture of two distributions provides a clear height threshold, while

the minimization of KS distance between observed and predicted data proves to be efficient for setting the minimum gap size. We use a Bayesian framework in which the model likelihood of each gap is expressed as a function of the unique environment local to the gap, highlighting the predominant role of the topographic exposure and waterlogging in determining gap size distribution. We expected that slope would also play an important role, with steeper slopes leading to larger gap sizes. However we found that a steeper slope lead to smaller gaps, as already highlighted by (Lobo and Dalling, 2014). We suggest that our modeling approach can be a basis for the development of large-scale methodologies using satellite data to understand gap phase dynamics at a regional scale, combining LiDAR and RaDAR remote sensing tools.

6 Supplementary information

Table 2. List of environmental variables, abbreviations, units, and values of the posteriors in univariate models for a height threshold equal to the 0.0001th percentile of the height distribution of the canopy.

Parameter	Abbreviation	Unit	Posterior value	Confidence interval (CI 95%)
Slope	Slope	°	0.119	[0.0416 ; 0.208]
Terrain Ruggedness Index	TRI	-	0.119	[0.083 ; 0.157]
TOPographic EXposure	TOPEX	-	-0.128	[-0.188 ; 0.00202]
Drained Area	DA	m ²	0.0843	[-0.0574 ; 0.179]
The Hydraulic Altitude	HAlt	m	-0.0135	[-0.04 ; 0.042]
HAND	HAND	-	-0.0615	[-0.152 ; 0.0162]

Table 3. List of environmental variables, abbreviations, units, and values of the posteriors in univariate models for a height threshold equal to the 0.001th percentile of the height distribution of the canopy.

Parameter	Abbreviation	Unit	Posterior value	Confidence interval (CI 95%)
Slope	Slope	°	0.0735	[-0.02 ; 0.15]
Terrain Ruggedness Index	TRI	-	0.0718	[0.04 ; 0.10]
TOPographic EXposure	TOPEX	-	-0.082	[-0.12 ; -0.05]
Drained Area	DA	m ²	-0.0176	[-0.09 ; 0.05]
The Hydraulic Altitude	HAlt	m	-0.0177	[-0.05 ; 0.02]
HAND	HAND	-	-0.003	[-0.08 ; 0.09]

Table 4. List of environmental variables, abbreviations, units, and values of the posteriors in univariate models for a height threshold equal to the 0.01th percentile of the height distribution of the canopy.

Parameter	Abbreviation	Unit	Posterior value	Confidence interval (CI 95%)
Slope	Slope	°	0.0975	[-0.02 ; 0.17]
Terrain Ruggedness Index	TRI	-	0.089	[0.05 ; 0.12]
TOPographic EXposure	TOPEX	-	-0.012	[-0.03 ; -0.32]
Drained Area	DA	m ²	-0.004	[-0.08 ; 0.05]
The Hydraulic Altitude	HAlt	m	0.063	[-0.04 ; 0.08]
HAND	HAND	-	-0.01	[-0.09 ; 0.06]

Acknowledgements. B.H. was supported by a grant from the Investing for the Future program (managed by the
360 350 French National Research Agency (ANR, labex CEBA, ref. ANR-10-LABX-0025).

References

- Asner, G. P. and Mascaro, J.: Mapping tropical forest carbon: Calibrating plot estimates to a simple LiDAR metric, *Remote Sensing of Environment*, 140, 614 – 624, doi:<http://dx.doi.org/10.1016/j.rse.2013.09.023>, <http://www.sciencedirect.com/science/article/pii/S003442571300360X>, 2014.
- 365 Asner, G. P., Kellner, J. R., Kennedy-Bowdoin, T., Knapp, D. E., Anderson, C., and Martin, R. E.: Forest canopy gap distributions in the southern Peruvian Amazon, *PloS one*, 8, e60 875, 2013.
- Bianchini, E., Garcia, C. C., Pimenta, J. A., and Torezan, J.: Slope variation and population structure of tree species from different ecological groups in South Brazil, *Anais da Academia Brasileira de Ciências*, 82, 643–652, 2010.
- 370 Bicknell, J. E., Phelps, S. P., Davies, R. G., Mann, D. J., Struebig, M. J., and Davies, Z. G.: Dung beetles as indicators for rapid impact assessments: evaluating best practice forestry in the neotropics, *Ecological Indicators*, 43, 154–161, 2014.
- Brokaw, N. V.: The definition of treefall gap and its effect on measures of forest dynamics, *Biotropica*, pp. 158–160, 1982.
- 375 Brokaw, N. V.: Gap-phase regeneration in a tropical forest, *Ecology*, 66, 682–687, 1985.
- Chapman, L.: Assessing topographic exposure, *Meteorological Applications*, 7, 335–340, 2000.
- Clauset, A., Shalizi, C. R., and Newman, M. E.: Power-law distributions in empirical data, *SIAM review*, 51, 661–703, 2009.
- Colson, F., GOND, V., FREYCON, V., Bogaert, J., and CEULEMANS, R.: Detecting natural canopy gaps in
380 Amazonian rainforest, *Bois et forêts des tropiques*, pp. 69–79, 2006.
- Denslow, J. S., Ellison, A. M., and Sanford, R. E.: Treefall gap size effects on above-and below-ground processes in a tropical wet forest, *Journal of Ecology*, 86, 597–609, 1998.
- Durrieu de Madron, L.: Mortalité des arbres en forêt primaire de Guyane française, *Bois et forêts des tropiques*, pp. 43–57, 1994.
- 385 Espírito-Santo, F. D., Keller, M. M., Linder, E., Oliveira Junior, R. C., Pereira, C., and Oliveira, C. G.: Gap formation and carbon cycling in the Brazilian Amazon: measurement using high-resolution optical remote sensing and studies in large forest plots, *Plant Ecology & Diversity*, 7, 305–318, 2014.
- Feeley, K. J., Davies, S. J., Ashton, P. S., Bunyavejchewin, S., Nur Supardi, M., Kassim, A. R., Tan, S., and Chave, J.: The role of gap phase processes in the biomass dynamics of tropical forests, *Proceedings of the Royal Society of London B: Biological Sciences*, 274, 2857–2864, doi:10.1098/rspb.2007.0954, <http://rspb.royalsocietypublishing.org/content/274/1627/2857>, 2007.
- 390 Ferry, B., Morneau, F., Bontemps, J.-D., Blanc, L., and Freycon, V.: Higher treefall rates on slopes and water-logged soils result in lower stand biomass and productivity in a tropical rain forest, *Journal of Ecology*, 98, 106–116, 2010.
- 395 Fisher, J. I., Hurr, G. C., Thomas, R. Q., and Chambers, J. Q.: Clustered disturbances lead to bias in large-scale estimates based on forest sample plots, *Ecology Letters*, 11, 554–563, 2008.
- Frolking, S., Palace, M. W., Clark, D., Chambers, J. Q., Shugart, H., and Hurr, G. C.: Forest disturbance and recovery: A general review in the context of spaceborne remote sensing of impacts on aboveground biomass and canopy structure, *Journal of Geophysical Research: Biogeosciences*, 114, 2009.

- 400 Gourlet-Fleury, S., Guehl, J.-M., and Laroussinie, O.: Ecology and management of a neotropical rainforest. Lessons drawn from Paracou, a long-term experimental research site in French Guiana, 2004.
- Guitet, S., Cornu, J.-F., Brunaux, O., Betbeder, J., Carozza, J.-M., and Richard-Hansen, C.: Landform and landscape mapping, French Guiana (South America), *Journal of Maps*, 9, 325–335, 2013.
- Hubbell, S. P., Foster, R. B., O'Brien, S. T., Harms, K., Condit, R., Wechsler, B., Wright, S. J., and De Lao, S. L.: Light-gap disturbances, recruitment limitation, and tree diversity in a neotropical forest, *Science*, 283, 554–557, 1999a.
- 405 Hubbell, S. P., Foster, R. B., O'Brien, S. T., Harms, K. E., Condit, R., Wechsler, B., Wright, S. J., and de Lao, S. L.: Light-Gap Disturbances, Recruitment Limitation, and Tree Diversity in a Neotropical Forest, *Science*, 283, 554–557, doi:10.1126/science.283.5401.554, <http://science.sciencemag.org/content/283/5401/554>, 1999b.
- 410 Kellner, J. R. and Asner, G. P.: Convergent structural responses of tropical forests to diverse disturbance regimes, *Ecology letters*, 12, 887–897, 2009.
- Kellner, J. R., Asner, G. P., Vitousek, P. M., Tweiten, M. A., Hotchkiss, S., and Chadwick, O. A.: Dependence of forest structure and dynamics on substrate age and ecosystem development, *Ecosystems*, 14, 1156–1167, 2011.
- 415 Korning, J. and Balslev, H.: Growth and mortality of trees in Amazonian tropical rain forest in Ecuador, *Journal of Vegetation Science*, 5, 77–86, 1994.
- Kuo, L. and Mallick, B.: Variable Selection for Regression Models, 1998.
- Lieberman, D., Lieberman, M., Peralta, R., and Hartshorn, G. S.: Mortality patterns and stand turnover rates in a wet tropical forest in Costa Rica, *The Journal of Ecology*, pp. 915–924, 1985.
- 420 Lloyd, J., Gloor, E. U., and Lewis, S. L.: Are the dynamics of tropical forests dominated by large and rare disturbance events?, *Ecology Letters*, 12, 2009.
- Lobo, E. and Dalling, J. W.: Spatial scale and sampling resolution affect measures of gap disturbance in a lowland tropical forest: implications for understanding forest regeneration and carbon storage, *Proceedings of the Royal Society of London B: Biological Sciences*, 281, 20133 218, 2014.
- 425 Mikita, T. and Klimánek, M.: Topographic Exposure and its Practical Applications, *Journal of Landscape Ecology*, 3, 42–51, 2012.
- Molino, J.-F. and Sabatier, D.: Tree Diversity in Tropical Rain Forests: A Validation of the Intermediate Disturbance Hypothesis, *Science*, 294, 1702–1704, doi:10.1126/science.1060284, <http://science.sciencemag.org/content/294/5547/1702>, 2001.
- 430 Negrón-Juárez, R. I., Chambers, J. Q., Hurr, G. C., Annane, B., Cocke, S., Powell, M., Stott, M., Goosem, S., Metcalfe, D. J., and Saatchi, S. S.: Remote Sensing Assessment of Forest Disturbance across Complex Mountainous Terrain: The Pattern and Severity of Impacts of Tropical Cyclone Yasi on Australian Rainforests, *Remote Sensing*, 6, 5633–5649, 2014.
- 435 Nelson, B. W., Kapos, V., Adams, J. B., Oliveira, W. J., and Braun, O. P.: Forest disturbance by large blowdowns in the Brazilian Amazon, *Ecology*, 75, 853–858, 1994.
- Nobre, A., Cuartas, L., Hodnett, M., Rennó, C., Rodrigues, G., Silveira, A., Waterloo, M., and Saleska, S.: Height above the nearest drainage—a hydrologically relevant new terrain model, *Journal of Hydrology*, 404, 13–29, 2011.

- 440 Puerta-Piñero, C., Muller-Landau, H. C., Calderón, O., and Wright, S. J.: Seed arrival in tropical forest tree fall gaps, *Ecology*, 94, 1552–1562, 2013.
- Riley, S. J.: Index That Quantifies Topographic Heterogeneity, 1999.
- Rüger, N., Huth, A., Hubbell, S. P., and Condit, R.: Response of recruitment to light availability across a tropical lowland rain forest community, *Journal of Ecology*, 97, 1360–1368, 2009.
- 445 Runkle, J. R.: Patterns of disturbance in some old-growth mesic forests of eastern North America, *Ecology*, pp. 1533–1546, 1982.
- SAGA, G.: System for automated geoscientific analyses, Available at: www.saga-gis.org/en/index.html, 2013.
- Schnitzer, S. A., van der Heijden, G., Mascaro, J., and Carson, W. P.: Lianas in gaps reduce carbon accumulation in a tropical forest, *Ecology*, 95, 3008–3017, 2014.
- 450 Team, R. C. et al.: R: A language and environment for statistical computing, 2013.
- Virkar, Y. and Clauset, A.: Power-law distributions in binned empirical data, *The Annals of Applied Statistics*, 8, 89–119, 2014.
- Wagner, F., Hérault, B., Stahl, C., Bonal, D., and Rossi, V.: Modeling water availability for trees in tropical forests, *Agricultural and Forest Meteorology*, 151, 1202–1213, 2011.
- 455 Whitmore, T.: Canopy gaps and the two major groups of forest trees, *Ecology*, 70, 536–538, 1989.
- Yee, T. W. et al.: The VGAM package for categorical data analysis, *Journal of Statistical Software*, 32, 1–34, 2010.