

# A Dirichlet process mixture model for Bayesian analysis of extreme avalanches

**Ophélie Guin<sup>1</sup>, Nicolas Eckert<sup>2</sup>, Liliane Bel<sup>3</sup>, Eric Parent<sup>3</sup>**

<sup>1</sup> Université de Lille

<sup>2</sup> UR ETNA, INRAE Grenoble,

<sup>3</sup> UMR MIA-Paris-Saclay, AgroParisTech, INRAE, université Paris-Saclay

AppliBUGS - 19 décembre 2023

## Extreme avalanches

Avalanches with **important runout distances / low stopping altitude**,  
⇒ provoke human and material losses.

Need methods to predict such avalanches.

Preliminary investigate :

- Change over time (e.g. impact of climate change),
- Stopping altitude distribution.

⇒ **Modeling stopping altitudes for extreme avalanches.**

## EPA dataset

Available data = **Enquête Permanente des Avalanches (EPA)**.

IC_014	Dept	Nom_Site	Height	IC_00001	Year	Intensite	Year_Start	IC_00002	Year_End	Observation	Year_observed_IC_000	
83641	70	SPARTE	1450	NO	30	1934	1935	1935	1935	2	FRAN RESSEC	Non observé
83642	70	SPARTE	1450	NO	30	1934	1935	1935	1935	2	FRAN RESSEC	Non observé
83643	70	SPARTE	1450	NO	30	1934	1935	1935	1935	2	FRAN RESSEC	Non observé
83644	70	SPARTE	1450	NO	30	1934	1935	1935	1935	2	FRAN RESSEC	Non observé
83645	70	SPARTE	1450	NO	30	1934	1935	1935	1935	2	FRAN RESSEC	Non observé
83646	70	SPARTE	1450	NO	30	1934	1935	1935	1935	2	FRAN RESSEC	Non observé
83647	70	SPARTE	1450	NO	30	1934	1935	1935	1935	2	FRAN RESSEC	Non observé
83648	70	SPARTE	1450	NO	30	1934	1935	1935	1935	2	FRAN RESSEC	Non observé
83649	70	SPARTE	1450	NO	30	1934	1935	1935	1935	2	FRAN RESSEC	Non observé
83650	70	SPARTE	1450	NO	30	1934	1935	1935	1935	2	FRAN RESSEC	Non observé
83651	70	SPARTE	1450	NO	30	1934	1935	1935	1935	2	FRAN RESSEC	Non observé
83652	70	SPARTE	1450	NO	30	1934	1935	1935	1935	2	FRAN RESSEC	Non observé
83653	70	SPARTE	1450	NO	30	1934	1935	1935	1935	2	FRAN RESSEC	Non observé
83654	70	SPARTE	1450	NO	30	1934	1935	1935	1935	2	FRAN RESSEC	Non observé
83655	70	SPARTE	1450	NO	30	1934	1935	1935	1935	2	FRAN RESSEC	Non observé
83656	70	SPARTE	1450	NO	30	1934	1935	1935	1935	2	FRAN RESSEC	Non observé
83657	70	SPARTE	1450	NO	30	1934	1935	1935	1935	2	FRAN RESSEC	Non observé
83658	70	SPARTE	1450	NO	30	1934	1935	1935	1935	2	FRAN RESSEC	Non observé
83659	70	SPARTE	1450	NO	30	1934	1935	1935	1935	2	FRAN RESSEC	Non observé
83660	70	SPARTE	1450	NO	30	1934	1935	1935	1935	2	FRAN RESSEC	Non observé
83661	70	SPARTE	1450	NO	30	1934	1935	1935	1935	2	FRAN RESSEC	Non observé
83662	70	SPARTE	1450	NO	30	1934	1935	1935	1935	2	FRAN RESSEC	Non observé
83663	70	SPARTE	1450	NO	30	1934	1935	1935	1935	2	FRAN RESSEC	Non observé
83664	70	SPARTE	1450	NO	30	1934	1935	1935	1935	2	FRAN RESSEC	Non observé
83665	70	SPARTE	1450	NO	30	1934	1935	1935	1935	2	FRAN RESSEC	Non observé
83666	70	SPARTE	1450	NO	30	1934	1935	1935	1935	2	FRAN RESSEC	Non observé
83667	70	SPARTE	1450	NO	30	1934	1935	1935	1935	2	FRAN RESSEC	Non observé
83668	70	SPARTE	1450	NO	30	1934	1935	1935	1935	2	FRAN RESSEC	Non observé
83669	70	SPARTE	1450	NO	30	1934	1935	1935	1935	2	FRAN RESSEC	Non observé
83670	70	SPARTE	1450	NO	30	1934	1935	1935	1935	2	FRAN RESSEC	Non observé
83671	70	SPARTE	1450	NO	30	1934	1935	1935	1935	2	FRAN RESSEC	Non observé
83672	70	SPARTE	1450	NO	30	1934	1935	1935	1935	2	FRAN RESSEC	Non observé
83673	70	SPARTE	1450	NO	30	1934	1935	1935	1935	2	FRAN RESSEC	Non observé
83674	70	SPARTE	1450	NO	30	1934	1935	1935	1935	2	FRAN RESSEC	Non observé
83675	70	SPARTE	1450	NO	30	1934	1935	1935	1935	2	FRAN RESSEC	Non observé
83676	70	SPARTE	1450	NO	30	1934	1935	1935	1935	2	FRAN RESSEC	Non observé
83677	70	SPARTE	1450	NO	30	1934	1935	1935	1935	2	FRAN RESSEC	Non observé
83678	70	SPARTE	1450	NO	30	1934	1935	1935	1935	2	FRAN RESSEC	Non observé
83679	70	SPARTE	1450	NO	30	1934	1935	1935	1935	2	FRAN RESSEC	Non observé
83680	70	SPARTE	1450	NO	30	1934	1935	1935	1935	2	FRAN RESSEC	Non observé
83681	70	SPARTE	1450	NO	30	1934	1935	1935	1935	2	FRAN RESSEC	Non observé
83682	70	SPARTE	1450	NO	30	1934	1935	1935	1935	2	FRAN RESSEC	Non observé
83683	70	SPARTE	1450	NO	30	1934	1935	1935	1935	2	FRAN RESSEC	Non observé
83684	70	SPARTE	1450	NO	30	1934	1935	1935	1935	2	FRAN RESSEC	Non observé
83685	70	SPARTE	1450	NO	30	1934	1935	1935	1935	2	FRAN RESSEC	Non observé
83686	70	SPARTE	1450	NO	30	1934	1935	1935	1935	2	FRAN RESSEC	Non observé
83687	70	SPARTE	1450	NO	30	1934	1935	1935	1935	2	FRAN RESSEC	Non observé
83688	70	SPARTE	1450	NO	30	1934	1935	1935	1935	2	FRAN RESSEC	Non observé
83689	70	SPARTE	1450	NO	30	1934	1935	1935	1935	2	FRAN RESSEC	Non observé
83690	70	SPARTE	1450	NO	30	1934	1935	1935	1935	2	FRAN RESSEC	Non observé

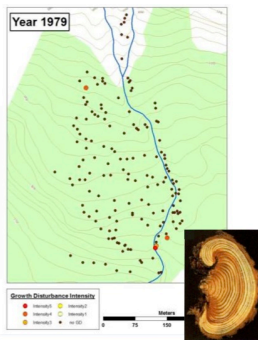
- Managed and developed by Inrae.
- Collection of data on avalanches (dates, snow cover, departure and arrival altitudes, type of avalanche, etc.).
- 3900 sites in 11 departments.
- More than 90,000 events available.

Extrait base EPA pour le site de Ressec

Time-limited database that only begins in 1900 + unclear at first.

## More data

How to retrieve more data for better estimations ?

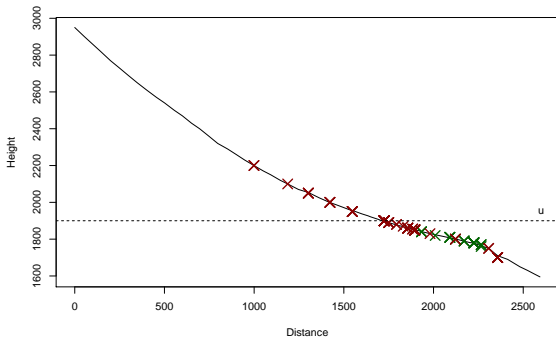


### Dendrochronological data.

- 1 Sample trees in an avalanche corridor,
- 2 detect impacted trees (in tree-rings) for each year,
- 3 determine avalanches years with spatial impact repartition,
- 4 set runout distance as the distance of the last impacted tree.

Difficulty : **censored data.**

## Data for Ressec site (Savoie).



**Goal** : probabilistic modeling of extreme avalanches as a function of time and stopping altitude.

Prerequisites :

- bivariate modeling,
- flexible modeling,
- taking into account that data may be censored,
- and non-stationary.

## Peak over threshold (POT) modeling for the stationary case (1)

Let  $X_1, X_2, \dots, X_n \stackrel{\text{iid}}{\sim} F$ ,  $X_j$  occur at regular time intervals.

$u$  a threshold level, we focus on the distribution  $F_u = \mathbb{P}(X - u \leq z | X > u)$ .

- $F_u$  converges in distribution to the Generalized Pareto's Distribution (GPD) :

$$\lim_{u \rightarrow u_\infty} F_u(z) = 1 - \left(1 + \xi \frac{z}{\sigma}\right)^{-1/\xi}.$$

[Pickands, 1975].

- Number of exceedances follows a Poisson's distribution.

## Peak over threshold (POT) modeling for the stationary case (2)

Relabel sub-sample  $(T_i, u + Z_i)$ , with  $i = 1, \dots, r$ , from ordered pairs  $(j, X_j)$  such that  $X_{T_j}$  exceeds  $u$ .

After re-normalization on  $\mathcal{A} = [0, 1] \times \mathbb{R}^+$ ,  $(T_i, Z_i)$  converges towards **homogeneous Poisson** process on  $\mathbb{R}^2$  with a separable intensity :

- homogeneous in time,
- GPD in the second component.

*[ Coles, 2001]*



## Extending for non-stationary case

Assumptions of the POT model :

- a sufficiently high threshold,
- independence for peaks over an acceptably high level.

⇒ Stationary behavior of the series.

Problem : stationary assumption in statistical avalanche analysis not realistic.

## Modeling

The point process  $(T_i, Z_i)$  are realizations such that :

- inhomogeneous in time,  
⇒ relax marginal Uniform distribution for a Beta distribution.
- frequency and severity of avalanches are not independent,  
⇒ in bayesian context components are independent *a priori* but not a *posteriori*.

**More flexibility** : take benefit from the classical POT model (with Beta x Pareto kernel) to generate weighted components of mixture model.

**Model interpretation** : study model uncertainty around the baseline of the classical POT approach for extremes.

## Intensity density of Poisson process

**Main question** : estimating the intensity density  $\lambda(t, z)$  of Poisson process.

Let the decomposition :

$$\lambda(\cdot) = \gamma f(\cdot),$$

$\gamma = \int_{\mathcal{A}} \lambda(t, z) dt dz$  the total intensity and  $f(\cdot)$  a density function.

Poisson process likelihood function :

$$L(\gamma, f(\cdot); \{(t_i, z_i) : i = 1, \dots, r\}) \propto \exp(-\gamma) \gamma^r \prod_{i=1}^r f(t_i, z_i),$$

**Estimation can be broken down into two independent problems** :

- $\gamma \rightarrow$  easy in Bayesian context.
- $f(\cdot) \rightarrow$  difficult.

## Kernel specification (1)

Specialists do not have much idea about density form in avalanches context.

⇒ Flexible model by relying on a nonparametric mixture :

$$f(t, z) = f(t, z; G) = \int_{\Theta} k(t, z|\theta) dG(\theta),$$

with  $k(t, z|\theta)$  a parametric density with parameter  $\theta$  and  $G$  a random mixing distribution.

Specification of bivariate kernel :

$$k(t, z|\theta) = k(t, z|\theta_1, \theta_2) = k_1(t|\theta_1)k_2(z|\theta_2),$$

where  $k_1$  and  $k_2$  independent before mixing.

## Kernel specification (2)

Beta distribution for kernel component over time :

$$k_1(t|\theta_1) = \frac{\Gamma(\tau)}{\Gamma(\tau\kappa)\Gamma(\tau(1-\kappa))} t^{\tau\kappa-1}(1-t)^{\tau(1-\kappa)-1},$$

with  $\kappa \in (0, 1)$  the mean and  $\tau > 0$  a scale parameter.

Generalized Pareto Distribution (GPD) for the exceedances :

$$k_2(z|\theta_2) = \frac{1}{\sigma} \left( 1 + \frac{\xi(z-u)}{\sigma} \right)^{-1/\xi-1}, \quad z \geq u,$$

with  $\theta_2 = (\sigma, \xi)$  with  $\sigma > 0$  and  $\xi > 0$ .

*Dendrochronological records considered as censored data.*

⇒  $k_2(z|\theta_2)$  will be replaced by  $\mathbb{P}(Z > z) = 1 - K_2(z|\theta_2)$  for these data.

## Bayesian Hierarchical Model

**Dirichlet process prior**  $DP(\alpha, G_0)$  is one of the most widely used Bayesian nonparametric priors.

[Ferguson, 1973]

$G_0$  the base distribution and  $\alpha$  controls how close the realization  $G$  is to  $G_0$ .

⇒ **Hierarchical model** :

$$\lambda(t, z) \equiv \lambda(t, z; G, \gamma) = \gamma f(t, z; G) = \gamma \int_{\Theta} k(t, z|\theta) dG(\theta)$$

$$G|\alpha, \theta \sim DP(\alpha, G_0),$$

## $\gamma$ inference

Marginal prior for  $\gamma$  :  $p(\gamma) \propto \gamma^{-1} \mathbf{1}_{\gamma > 0}$ .

[Kottas and Behseta, 2010].

$\Rightarrow$  Proper posterior distribution  $p(\gamma|t, z)$  is a gamma( $n, 1$ ).

## Dirichlet process algorithm (1)

Dirichlet process realizations are discrete with probability one,

⇒ model can be viewed as infinite mixtures [Ferguson, 1983].

Equivalent model obtained by taking the limit as  $L$  goes to infinity of finite mixture models with  $L$  components :

$$(t_i, z_i) | \kappa_{L_i}, \tau_{L_i}, \sigma_{L_i}, \xi_{L_i} \sim k_1(t_i | \kappa_{L_i}, \tau_{L_i}) k_2(z_i | \sigma_{L_i}, \xi_{L_i}), i = 1, \dots, r$$

$$L_i | \mathbf{p} \sim \text{Discrete}(p_1, \dots, p_L)$$

$$\mathbf{p} | \alpha \sim \text{Dirichlet}(\alpha/L, \dots, \alpha/L)$$

$$\theta_l = (\kappa_l, \tau_l, \sigma_l, \xi_l) \sim \mathbf{G}_0(\theta_l | \psi), l = 1, \dots, L,$$

$L_i$  represent à “latent class” associated with observation  $(t_i, y_i)$ .



## Dirichlet process algorithm (2)

Posterior inference for this type of model based on the Chinese Restaurant Process sampler [Neal, 2000].

Here  $G_0$  base measure is a non-conjugate prior for  $\theta$ ,  
⇒ more difficulties and use of numerical techniques.

Algorithm used for inference is Algorithm 8 of [Neal, 2000].

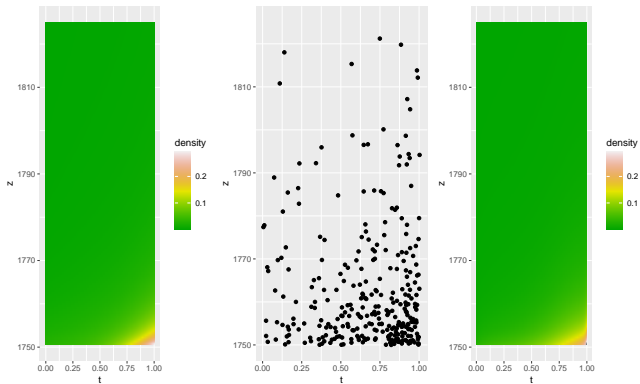
## Base distribution

The different base distribution components are *a priori* independent :

$$G_0(\theta) = G_0(\kappa, \tau, \sigma, \xi) = G_0^\kappa(\kappa)G_0^\tau(\tau)G_0^\sigma(\sigma)G_0^\xi(\xi)$$

- $G_0^\kappa$  is an uniform distribution,
- $G_0^\tau$  and  $G_0^\sigma$  are inverse-gamma distributions with fixed shape parameter,
- $G_0^\xi$  is an exponential distribution.

## Simulations



## Datasets

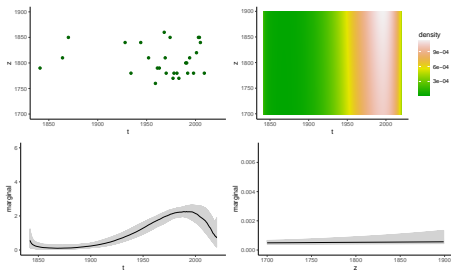
Two types of data for the Ressec site :

- 67 EPA data between 1901 and 2022,
- 28 dendrochronological data from 1840.

**Data pre-treatment** : select only the avalanche from the EPA base (more precise) when avalanches appear in the two bases.

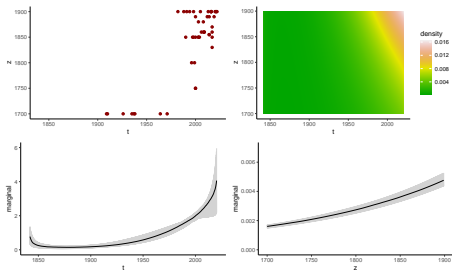
Threshold  $u = 1900$ .

## Model apply to dendrochronological data



- Increase in avalanches over time.
- Unclear results for stopping altitudes.
  - ⚠ Only worked with censored data.

## Model apply to EPA data

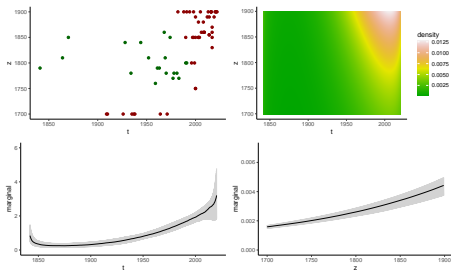


- Probability of extreme avalanches increases with time (with acceleration in the 2000s).

⚠ Few data before the 1980s.

- Stopping altitudes probability increases for higher altitude.

## Model apply to EPA and dendrochronological data



- Dendrochronological data provide information (before 1980 and for stopping altitudes greater than 1700m).
- Probability of extreme avalanches increases with time.
- Stopping altitudes probability increases but acceleration in the 1950s.

## Conclusion and perspectives

### Non-parametric modeling of extreme non-stationary and censored data.

Perspectives :

- 1 Here, restriction to the case  $\xi > 0$ ,  
⇒ extend to all values of  $\xi$ .
- 2 Return period estimations.



## References



Coles, S. (2001)

An Introduction to Statistical Modeling of Extreme Values

*Springer*



Ferguson, T. S. (1973)

A bayesian analysis of some nonparametric problems.

*The Annals of Statistics, 1 :209–230.*



Neal, R. (2000).

Markov chain sampling methods for dirichlet process mixture models.

*Journal of Computational and Graphical Statistics, 9(2) :249–265.*



Pickands, J. (1975).

Statistical inference using extreme order statistics.

*The Annals of Statistics, 3(1) :119–131.*