

Modeling extreme values of processes observed at irregular  
time steps.  
Application to significant wave height

Pierre Ailliot

Laboratoire de Mathématiques de Bretagne Atlantique  
Université de Brest

Joint work with Nicolas Raillard (IFREMER)

# Plan of the talk

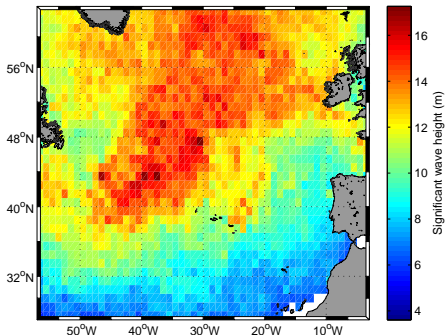
- 1 Introduction
- 2 Method
  - Usual methods in the iid case
  - Generalization to dependent sequences
- 3 Analysis of Hs data
- 4 Conclusion

# Outline

- 1 Introduction
- 2 Method
  - Usual methods in the iid case
  - Generalization to dependent sequences
- 3 Analysis of Hs data
- 4 Conclusion

## Introduction

- Hs : **significant** wave height, parameter related to wave energy in a sea state
  - Historical definition: mean height of the one-third highest waves
  - Current definition: four times the standard deviation of the sea surface elevation
- Aim: develop a method to estimate the extremal properties of Hs
  - Quantities of interest: **return levels**, storm durations,...
  - Example: 20-year return levels in the North Atlantic (method explained hereafter)



# Introduction

- What are the data available to estimate the extremal properties of Hs?
  - **Reanalysis data**
    - Rich space-time sampling
    - Standard statistical method apply!
    - Underestimation for high Hs?
  - **Buoy data**
    - Sparse spatial sampling
    - Rich time sampling
    - Most reliable for high Hs?
  - **Satellite altimeter data**
    - Sparse space-time sampling
    - Good spatial coverage
    - Quality of the data for high Hs?

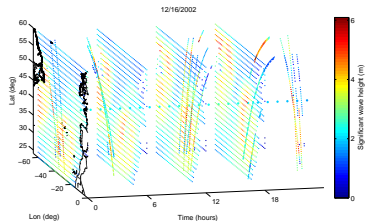


Figure: Data available on 15/12/2002.

# Hs data

- Single site analysis

- Buoy Brittany (47.5 N, 8.5 W)
  - 10 years of data available,  $\Delta t = 1h$
  - 7.7% of missing data
  - Breakdowns linked to extreme events?
- Reanalysis data: ERA interim (ECMWF)
  - Same years than buoy,  $\Delta t = 6h$
- Satellite data
  - 15 years of data available, 7 different satellites
  - Closest observations in the satellite tracks which intersect a  $3^{\circ} \times 3^{\circ}$  box

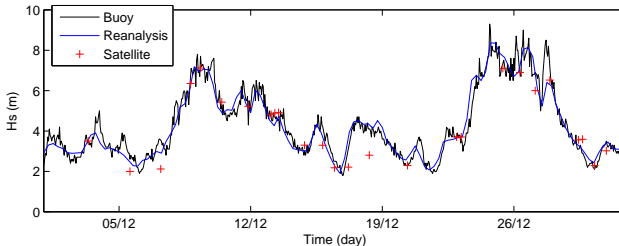


Figure: Hs data for dec. 1999

# Hs data

- Single site analysis

- Buoy Brittany (47.5 N, 8.5 W)
  - 10 years of data available,  $\Delta t = 1h$
  - 7.7% of missing data
  - Breakdowns linked to extreme events?
- Reanalysis data: ERA interim (ECMWF)
  - Same years than buoy,  $\Delta t = 6h$
- Satellite data
  - 15 years of data available, 7 different satellites
  - Closest observations in the satellite tracks which intersect a  $3^{\circ} \times 3^{\circ}$  box



- How can we analyze the extremal behavior of such time series?
- Need for methods which can deal with **missing data** and **irregular time sampling**
- Can we get useful information on extreme Hs from satellite data?

# Outline

- 1 Introduction
- 2 **Method**
  - Usual methods in the iid case
  - Generalization to dependent sequences
- 3 Analysis of Hs data
- 4 Conclusion

## Probabilistic background (iid case)

- Let  $M_n = \max_{i=1, \dots, n} X_i$  with  $(X_j)$  i.i.d. sample with c.d.f.  $G$  ;

## Probabilistic background (iid case)

- Let  $M_n = \max_{i=1, \dots, n} X_i$  with  $(X_i)$  i.i.d. sample with c.d.f.  $G$  ;
- If  $\mathbb{P}\left\{\frac{M_n - b_n}{a_n} \leq x\right\} = G^n(a_n x + b_n) \rightarrow F(x)$ , then  $F$  is a max-stable distribution ie

$$F^n(\alpha_n x + \beta_n) = F(x)$$

# Probabilistic background (iid case)

- Let  $M_n = \max_{i=1, \dots, n} X_i$  with  $(X_i)$  i.i.d. sample with c.d.f.  $G$  ;
- If  $\mathbb{P}\left\{\frac{M_n - b_n}{a_n} \leq x\right\} = G^n(a_n x + b_n) \rightarrow F(x)$ , then  $F$  is a max-stable distribution ie

$$F^n(\alpha_n x + \beta_n) = F(x)$$

- Fisher-Tippett (1928), Gnedenko (1943): Max-stable distributions have cdf

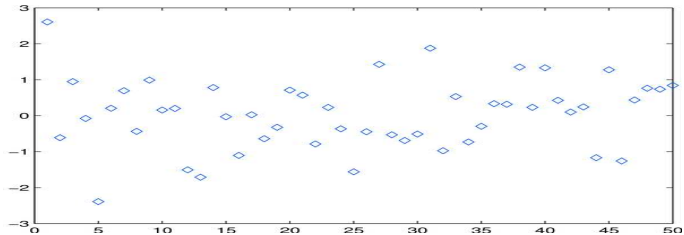
$$F(x; \mu, \sigma, \xi) = \exp \left[ - \left( 1 + \xi \frac{x - \mu}{\sigma} \right)_+^{-1/\xi} \right] \in \text{GEV}(\mu, \sigma, \xi)$$

This **GEV** distribution includes the **Weibull** (finite upper bound,  $\xi < 0$ ), **Fréchet** (heavy tail,  $\xi > 0$ ) and **Gumbel** ( $\xi = 0$ ) distributions

- The GEV is a natural distribution for modeling the maximum of a large number of i.i.d. random variables

# Usual methods for analyzing extremes of i.i.d. sequences

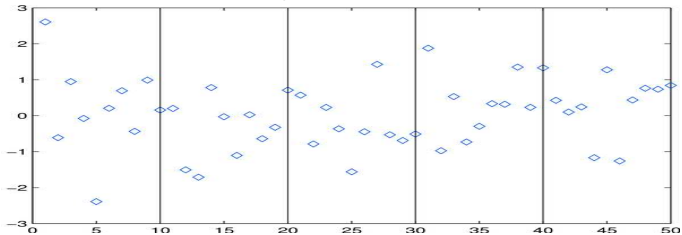
- Block maxima method



# Usual methods for analyzing extremes of i.i.d. sequences

- **Block maxima method**

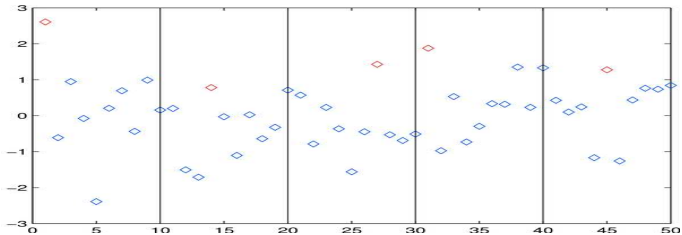
- Group the data into blocks of equal lengths (usually one year to remove seasonal effects)



# Usual methods for analyzing extremes of i.i.d. sequences

## • Block maxima method

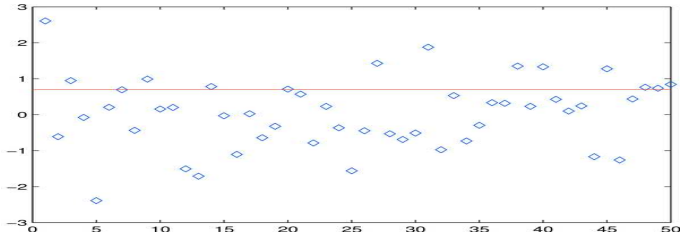
- Group the data into blocks of equal lengths (usually one year to remove seasonal effects)
- Fit a GEV distribution to the sample of block maxima



- Waste of data (only the yearly maxima are kept to fit the GEV)
  - Only 10 observations to fit the GEV for buoy data although more storms are observed!

# Usual methods for analyzing extremes of i.i.d. sequences

- **Block maxima method**
  - Group the data into blocks of equal lengths (usually one year to remove seasonal effects)
  - Fit a GEV distribution to the sample of block maxima
  - Waste of data (only the yearly maxima are kept to fit the GEV)
    - Only 10 observations to fit the GEV for buoy data although more storms are observed!
- **Peaks Over Threshold (POT)**
  - Choose a high threshold  $u$



# Usual methods for analyzing extremes of i.i.d. sequences

## • Block maxima method

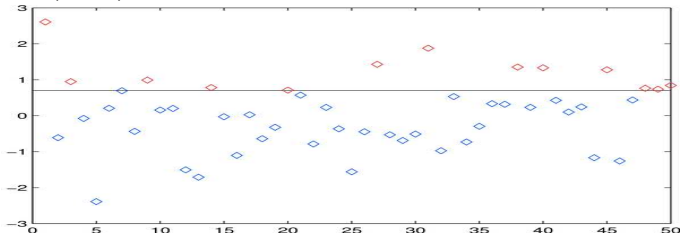
- Group the data into blocks of equal lengths (usually one year to remove seasonal effects)
- Fit a GEV distribution to the sample of block maxima
- Waste of data (only the yearly maxima are kept to fit the GEV)
  - Only 10 observations to fit the GEV for buoy data although more storms are observed!

## • Peaks Over Threshold (POT)

- Choose a high threshold  $u$
- Fit a GPD distribution to the exceedances above  $u$

$$X_i - u | X_i > u \sim_{i.i.d.} \text{GPD}(\tilde{\sigma}, \xi)$$

with  $\lambda = \mathbb{P}(X_i > u)$  an extra unknown parameter



# Usual methods for analyzing extremes of i.i.d. sequences

## • Block maxima method

- Group the data into blocks of equal lengths (usually one year to remove seasonal effects)
- Fit a GEV distribution to the sample of block maxima
- Waste of data (only the yearly maxima are kept to fit the GEV)
  - Only 10 observations to fit the GEV for buoy data although more storms are observed!

## • Peaks Over Threshold (POT)

- Choose a high threshold  $u$
- Fit a GPD distribution to the exceedances above  $u$

$$X_i - u | X_i > u \sim_{i.i.d} \text{GPD}(\tilde{\sigma}, \xi)$$

with  $\lambda = \mathbb{P}(X_i > u)$  an extra unknown parameter

- Similar results when fitting a **censored GEV distribution**

$$u\mathbf{1}[X_i \leq u] + X_i\mathbf{1}[X_i > u] = u\mathbf{1}[Y_i \leq u] + Y_i\mathbf{1}[Y_i > u] \text{ with } Y_i \sim_{i.i.d} \text{GEV}(\mu, \sigma, \xi)$$

- GEV and GPD distributions lead to similar tail approximations. If  $x \rightarrow x+$  then

$$\exp \left[ - \left( 1 + \xi \frac{x - \mu}{\sigma} \right)_+^{-1/\xi} \right] \approx 1 - \left( 1 + \xi \frac{x - \mu}{\sigma} \right)_+^{-1/\xi}$$

- Successive exceedances are generally dependent

- A "declustering" step is generally applied to keep only one value per "storm"
- Difficulties: arbitrary rules to define the "storms", waste of data...
- Alternative...model  $\{Y_i\}$  as a **max-stable process** and keep all exceedances

# Max-stable processes

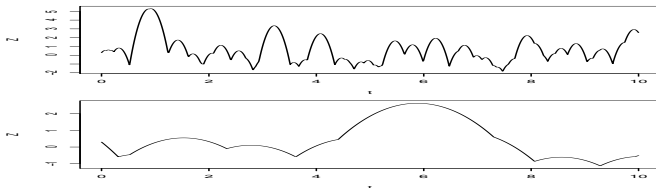
- Max-stable processes are natural approximations of processes over high thresholds
  - Similar to the GEV in the i.i.d. case
  - No unified parametric representation for max-stable processes
- A particular max-stable process: the **Gaussian extreme value process**
  - Assume that  $\{Y_t\}$  has GEV marginal distribution with parameter  $(\mu, \sigma, \xi)$ .
  - Let

$$Z_t = -\frac{1}{\log(F(Y_t; \mu, \sigma, \xi))},$$

$\{Z_t\}$  has unit Fréchet marginal distribution. Assume that

$$Z_t = \max \left\{ \frac{\zeta_i}{\nu\sqrt{2\pi}} \exp\left(-\frac{(s_i - t)^2}{2\nu^2}\right) \right\},$$

with  $\{(\zeta_i, s_i), i \geq 1\}$  the points of a Poisson process with intensity  $\zeta^{-2}d\zeta \times ds$ .



- $\nu \rightarrow 0$  [resp.  $\nu \rightarrow +\infty$ ] corresponds to independence [resp. perfect dependence]

# Max-stable processes

- Max-stable processes are natural approximations of processes over high thresholds
  - Similar to the GEV in the i.i.d. case
  - No unified parametric representation for max-stable processes
- A particular max-stable process: the **Gaussian extreme value process**
  - Assume that  $\{Y_t\}$  has GEV marginal distribution with parameter  $(\mu, \sigma, \xi)$ .
  - Let

$$Z_t = -\frac{1}{\log(F(Y_t; \mu, \sigma, \xi))},$$

$\{Z_t\}$  has unit Fréchet marginal distribution. Assume that

$$Z_t = \max \left\{ \frac{\zeta_i}{\nu\sqrt{2\pi}} \exp \left( -\frac{(s_i - t)^2}{2\nu^2} \right) \right\},$$

with  $\{(\zeta_i, s_i), i \geq 1\}$  the points of a Poisson process with intensity  $\zeta^{-2}d\zeta \times ds$ .

- The full joint distribution is not tractable but **bivariate** distributions are

$$P(Y_{t_1} \leq y_{t_1}, Y_{t_2} \leq y_{t_2}) = \exp \left[ -\frac{1}{z_{t_1}} \Phi \left( \frac{a}{2} + \frac{1}{a} \log \frac{z_{t_2}}{z_{t_1}} \right) - \frac{1}{z_{t_2}} \Phi \left( \frac{a}{2} + \frac{1}{a} \log \frac{z_{t_1}}{z_{t_2}} \right) \right]$$

with  $a = \frac{|t_1 - t_2|}{\nu}$ ,  $\Phi$  the cdf of the standard normal distribution and  $z_{t_i} = \frac{-1}{\log F(Y_{t_i}; \mu, \sigma, \xi)}$ .

- **Continuous time process**...permits to deal with irregular time sampling

## Proposed methodology

- Let  $\{X_t\}_{t \in \{t_1, \dots, t_n\}}$  be a stationary process observed at time  $(t_1, \dots, t_n)$
- Chose a high threshold  $u$
- Model the process censored at  $u$  as a **censored max-stable process**

$$u\mathbf{1}[X_{t_i} \leq u] + X_{t_i}\mathbf{1}[X_{t_i} > u] = u\mathbf{1}[Y_{t_i} \leq u] + Y_{t_i}\mathbf{1}[Y_{t_i} > u]$$

with  $\{Y_t\}$  a **Gaussian extreme value process** with parameter  $\theta = (\mu, \sigma, \xi, \nu)$

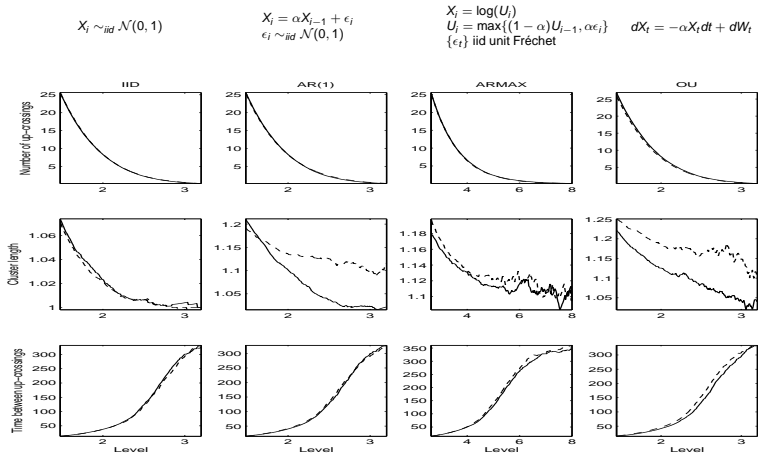
- Process with marginal distribution  $GEV(\mu, \sigma, \xi)$ ,  $\nu$  describes the time structure
- Estimate  $\theta$  by maximizing the **pairwise likelihood** function

$$PL(\theta; \tilde{y}_{t_1}, \dots, \tilde{y}_{t_n}) = \prod_{i=1}^{n-1} p_{\tilde{y}}(\tilde{y}_{t_i}, \tilde{y}_{t_{i+1}}; \theta)$$

of the censored sequence  $(\tilde{y}_{t_1}, \dots, \tilde{y}_{t_n})$  with  $\tilde{y}_t = u\mathbf{1}[y_t \leq u] + y_t\mathbf{1}[y_t > u]$

- Consistent estimator
- Better results when using only closest neighbors in the pairwise likelihood function
- Use simulations of the fitted model to estimate quantities of interest
  - Distribution of the estimator (parametric bootstrap)
  - Return values, length of the storms, number of storm per year,...

# Validation on classical time series models



**Figure:** Comparison of the extremal behavior of classical processes (solid line) against the fitted Gaussian extreme value process (dashed line). Results obtained by simulating 1000 years of each model (one observation per day).

# Outline

- 1 Introduction
- 2 Method
  - Usual methods in the iid case
  - Generalization to dependent sequences
- 3 Analysis of Hs data
- 4 Conclusion

# Data

- Focus on December to remove seasonality

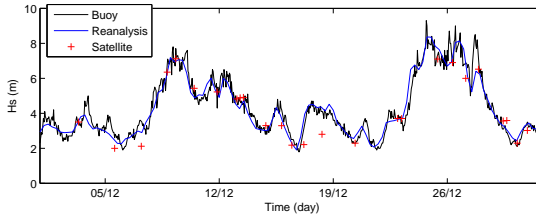
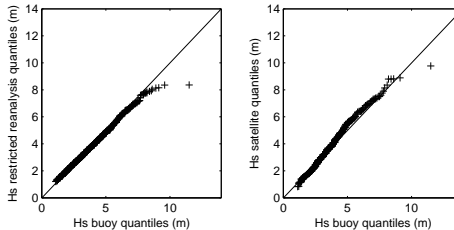


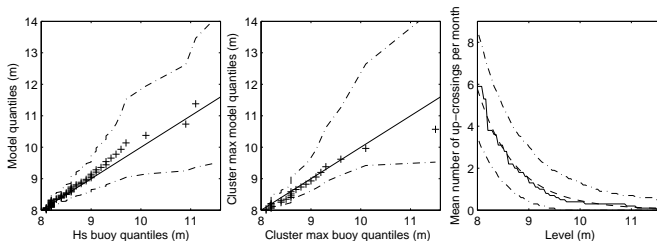
Figure: Hs data for dec. 1999

- The quantiles of the satellite data are the closest to the ones of the buoy



## A few results

- The model has been fitted to the different data sets
- It seems to provide a realistic extrapolation of the extremal properties of the data
  - Comparison of statistics for the buoy data and the fitted model



## A few results

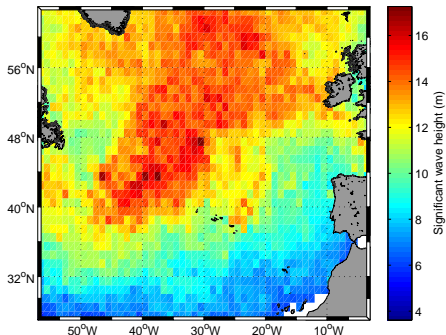
- The model has been fitted to the different data sets
- It seems to provide a realistic extrapolation of the extremal properties of the data
- Comparison of parameters and return values for the different data sets

	Reanalysis	Buoy	Satellite
Threshold			
$u$ (m)	6	8	6
Nb obs $> u$	74	59	48
Parameter values			
$\mu$	2.19 [-6.85,3.80]	5.12 [-7.3,6.67]	3.80 [2.23,4.47]
$\sigma$	1.70 [0.81,8.83]	0.50 [0.11,7.23]	1.33 [0.86,2.69]
$\xi$	-0.17 [-0.6,-0.02]	0.07 [-0.40,0.32]	0.01 [-0.25,0.17]
$\nu$	0.12 [0.07,0.18]	1e-3 [8e-4,2e-3]	0.05 [0.03,0.08]
Return levels			
$q_{10}$	8.2[7.3,9.1]	10.4[9.3,12.2]	12.1[10.7,14.6]
$q_{20}$	8.5[7.5,9.6]	10.8[9.4,13.8]	12.9[11.1,16.4]
$q_{50}$	8.8[7.6,10.1]	11.4[9.6,16.0]	14.0[11.5,19.1]
$q_{100}$	9.0[7.7,10.6]	11.8[9.7,18.8]	14.7[11.7,21.6]

- Confidence intervals (computed using parametric bootstrap) are wide!
- Gumbel type distributions ( $\xi = 0$ ) except for the restricted reanalysis data ( $\xi < 0$ )
  - Lower tail for the reanalysis data?
- Higher value of  $\nu$  for the reanalysis data (smoother, longer storms)
- Buoy and satellite lead to similar fitted models

## A few results

- The model has been fitted to the different data sets
- It seems to provide a realistic extrapolation of the extremal properties of the data
- Comparison of parameters and return values for the different data sets
- 20 years return levels in the North Atlantic computed using satellite data
  - Independent analysis at each location, no spatial information





# Outline

- 1 Introduction
- 2 Method
  - Usual methods in the iid case
  - Generalization to dependent sequences
- 3 Analysis of Hs data
- 4 Conclusion

## Conclusion

- The proposed methodology permits to analyze the extremal behavior of **dependent** sequences
- Still valid in presence of **missing data** or **irregular time sampling**
- Theoretical results show that the estimator is consistent
- The method has been successfully validated on classical time series models
- Similar models are identified on satellite and buoy data
- Longer storm and lighter tails are identified on reanalysis data
- Preprint available on my web page  
Raillard N., Ailliot P., Yao J.F., Modeling extreme values of processes observed at irregular time steps. Application to significant wave height. *To appear in the Annals of Applied Statistics.*

-  **N. Raillard, P. Ailliot, and J.F. Yao.**  
Modelling extreme values of processes observed at irregular time step. application to significant wave height.  
*Submitted, 2011.*
-  **R. L. Smith.**  
Max-stable processes and spatial extremes.  
*Unpublished, 1990.*