

UNIVERSITE CATHOLIQUE DE LOUVAIN

FACULTE DES SCIENCES

ECOLE DE STATISTIQUE, BIOSTATISTIQUE
ET SCIENCES ACTUARIELLES

The text is really very long... Please look
for ways to reduce. "Less is more."



TEMPORAL ANALYSIS OF THE EVOLUTION OF EXTREME VALUES USING
CLIMATOLOGICAL DATA

Promoteur : Johan SEGERS

□ Lecteurs : Anna KIRILIOUK
Michel CRUCIFIX

Mémoire présenté en vue de l'obtention du

Master en statistiques, orientation générale

par : **Antoine Pissoort**

Juin 2017

Abstract

This thesis aims to analyse extreme temperatures from Uccle and assess their nonstationarity. trend in the location parameter of the temperatures assessing the climate warming. After having proven it by hand of splines's derivatives that a correction for simultaneous intervals led to non-significant changes, we will do it by hand of the Extreme Value Theory (EVT) that there is indeed an upward trend in the location parameter of the temperatures assessing the climate warming. First of all, we will do introductory analysis of the trend... and we will discover that.. The analysis will focus on yearly maxima and hence we will go through with EVT by defining and presenting the usual methods such as GEV.

Regarding the computations, we took advantage of a high-level language (c++) to make our analysis efficient and also made use of parallel computing to decrease computation time for time consuming...

Keywords • Extreme Value Theory • block-maxima model • peaks-over-threshold method • trend analysis • Bayesian inference • Generalized Additive Models with splines smoothing • Neural Networks • nonstationary models • Markov Chain Monte Carlo • Hamiltonian Monte Carlo • Parallel computing • R package • Shiny application

Acknowledgements

I would first like to thank my thesis supervisor Johan Segers for all his help and his guidance during this whole year. The repeated appointments we have had

I also would like to thank the "Institut Royal de Météorologie" (IRM) of Belgium for his help and his guidance but also for his provided quality datasets.

Finally, I want to thank my family and my friends, but also Bernadette for her support and all the time I have spent writing in her room for my thesis but also during my whole academic studies.

Contents

Introduction	x
I Theoretical Framework : Extreme Value Theory	3
1 Method of Block-Maxima	4
1.1 Preliminaries	5
1.2 Extremal Types Theorem	7
1.3 Characterization of the GEV distributions : 3 Types	8
1.4 Applications : Examples of Convergence to GEV	10
1.5 Maximum Domain of Attraction	13
1.5.1 Domain of attraction for the 3 types of GEV	14
1.5.2 Closeness under tail equivalence property	17
1.5.3 Domain of attraction of the GEV	18
1.6 The Concepts of Return Levels and Return Periods	18
1.7 Inference	19
1.7.1 Likelihood-based Methods	19
1.7.2 Other Methods	21
1.8 Model Diagnostics : Goodness-of-Fit	21
1.8.1 Return Level Plot	21
2 Peaks-Over-Threshold Methods	24
2.1 Preliminaries: Intuitions	25
2.2 Characterization of the Generalized Pareto Distribution	26
2.2.1 Outline proof of the GPD and justification from GEV	27
2.2.2 Dependence of the scale parameter σ	28
2.2.3 Three different types of GPD and duality with GEV	28
2.2.4 Examples of the GPD as limiting distribution for exceedances	29
2.3 Return Levels	29
2.4 Point Process Approach	30

2.4.1	Non-homogeneous Poisson Process	30
2.5	Inference	30
2.5.1	Likelihood-based Methods	31
2.5.2	Profile Likelihood	31
2.5.3	Other Methods	31
2.5.4	Estimators Based on Extreme Order Statistics for EVI	31
2.5.5	The Probability-Weighted-Moment Estimator	32
2.5.6	Estimators based on Generalized Quantile	33
2.6	Threshold Selection (Methods)	33
2.6.1	Standard Threshold choice for the excess models	33
2.6.2	"Varying" Threshold : Mixture Models	36
2.6.2.1	Nonstationary extremes	36
3	Relaxing The Independence Assumption	38
3.1	Stationary Extremes	38
3.1.1	The extremal index	40
	Clusters of exceedances	40
	New parameters	40
	Return levels	41
3.1.2	Tail dependence	41
3.1.3	Modelling : Threshold Models	41
3.1.4	Applications	42
3.2	Non-Stationary Extremes	42
3.2.1	Block-Maxima	43
3.2.2	Diagnostics	43
3.3	Model Comparisons	43
3.3.1	Statistical Tools	43
3.4	Return Levels	43
3.5	Inference	44
4	Neural Network (and others)	45
4.1	Improvements For Modelling Non-stationary Sequences	45
4.1.1	Generalized Likelihood Methods	45
4.2	Neural-Network Based Inference	45
4.3	Bagging	47
4.4	Bootstrap Methods in EVT	47
4.4.1	Moving Block Bootstrap	48

4.5	Markov models	48
5	Bayesian Methods	49
5.1	Prior Elicitation	50
5.1.1	Non-informative Priors	50
5.1.2	Informative Priors	51
5.2	Bayesian Computation : Markov Chains	52
5.2.1	Algorithms	52
5.2.2	Hamiltonian Monte Carlo	53
5.2.3	Computational efficiency comparison	53
5.3	Convergence Diagnostics	53
5.3.1	Proposal Distribution	53
5.3.2	The problem of auto and cross-correlations in the chains	54
5.4	Posterior Predictive	54
5.5	Bayesian Predictive Accuracy for Model Validation	55
5.5.1	Cross-validation for predictive accuracy	55
5.6	Bayesian Inference ?	57
5.6.1	Bayesian Credible Intervals	57
5.6.2	Distribution of Quantiles : Return Levels	57
5.7	Bayesian Model Averaging	57
II	Experimental Framework : Nonstationary Extreme Value Analysis of Maximum Temperatures	58
6	Introduction to the Analysis	59
	Repository for the code : R Package	59
	Visualization Tool : Shiny Application	59
6.1	Presentation of the Analysis : Temperatures from Uccle	60
6.1.1	Open shelter vs Closed shelter	60
6.1.2	Comparisons with freely available data	60
6.2	First Analysis : Block-Maxima	61
6.2.1	Descriptive Analysis	61
6.2.2	First visualization with simple models	61
6.2.3	Deeper Trend Analysis : Splines derivatives in GAM	62
	Pointwise vs Simultaneous intervals	62
	Methodology	62
	Final Results	63

6.3	Comments and Structure of the Analysis	64
7	First Analysis by GEV	66
	R packages for EVT	66
7.1	Estimation of the Model	67
7.1.1	Maximum Likelihood	67
7.1.2	Other Methods	68
7.2	Diagnostics	68
7.3	Return Levels	68
7.3.1	Profile Likelihood	68
7.4	Comments and Comparisons with POT	68
8	Stationary and Nonstationary GEV Analysis	69
8.1	Stationary Analysis	69
8.2	Nonstationary Analysis	69
8.3	Improvements with Neural Networks	69
9	Bayesian Analysis	70
9.1	From evdbayes R package : MH algorithm	70
9.2	From Our Functions (R package)	70
9.3	From HMC algorithm using STAN language	70
9.4	Comparisons	71
9.4.1	STAN	71
9.5	Comparison with frequentists results	71
10	Conclusion	72
	Appendix	73
A	Statistical tools for Extreme Value Theory	75
A.1	Tails of the distributions	75
A.2	Convergence concepts	76
A.3	Varying functions	77
A.4	Diagnostic Plots : Quantile and Probability Plots	77
B	Bayesian Methods	79
B.1	Algorithms	79
B.1.1	Metropolis–Hastings Algorithm	79

B.1.2	Gibbs Sampler	80
B.1.3	Hamiltonian Monte Carlo	81
C	Other Figures and Tables	84
C.1	GEV : Influence of the Parameters on the shape of the distribution	84
C.2	Introduction of the Practical Analysis (section 6)	85
C.3	Analysis by GEV	85
D	Github Repository Structure	88

List of Figures

1.1	GEV distribution with the normal as benchmark (dotted lines) and a zoom on the parts of interest to better visualize the behaviour in the tails. In green, we retrieve the Gumbel distribution ($\xi = 0$). In red, we retrieve the Weibull-type ($\xi < 0$) while in blue, we get the Fréchet-type ($\xi > 0$). The endpoints for the Weibull and the Fréchet are denoted by the red and blue filled circles respectively.	10
4.1	<i>Neural Network applied to GEV. Figure made with <code>tikzpicture</code> and based on Cannon (2010)</i>	46
6.1	representing the yearly maxima together with three first models trying to represent the trend. Note that shaded grey line representing the standard errors (and not a confidence interval) of the linear trend.	61
6.2	displays draws from the posterior distribution of the model. Notice the large uncertainty associated to the posterior draws (grey lines). The coverages are calculated for $M = 10^5$ simulations.	63
6.3	Plots of the first derivative $f'(\cdot)$ of the estimated splines on the retained GAM model. Grey area represents the 95% confidence interval. Sections of the spline where the confidence interval does not include zero are indicated by thicker sections.	64
7.1	68
7.2	The left and right vertical dotted lines represent respectively the minimum and the maximum value of the yearly maxima series.	68
C.1	GEV distribution for different values of the three parameters	84
C.2	ACF and PACF for the residuals of the fitted GAM model with assumed independent errors	85
C.3	Diagnostics of the chosen GAM model with MA(1) errors, based on the residuals. . . .	86
C.4	Series of annual maxima together with the fitted GAM model (in green). Thicker lines indicate that the increase is significant for <u>pointwise</u> confidence interval. Shaded area represent the "95%" interval which looks quite narrow.	86
C.5	The two horizontal dotted lines represent the 95% (above) and 99% (below) confidence intervals when we take the intersection on the horizontal axis.	87

List of Abbreviations

For convenience, we place a list of all the abbreviations we will use in the text. However, these will always be defined in their first occurrence into the text.

df	distribution function
EVI	Extreme Value Index (ξ)
EVT	Extreme Value Theory
GEV	Generalized Extreme Value
GPD	Generalized Pareto Distribution (function)
MCMC	Marko Chain Monte Carlo
MH	Metropolis-Hastings (algorithm)
R.V.	Random Variable
TN	Temperature miNimum
TX	Temperature maXimum

Introduction

Unlike his counterparts (see for example credit risk analysis, financial applications,...), the Extreme Value Theory (EVT) applied on the broad environmental area like here for the meteorological data, has strong impacts on the people lives

An important question is still whether climate changes caused by anthropogenic activities will change the intensity and frequency of extreme events [Milly et al. \(2008a\)](#).

The problem we are here facing in climate change evidence is that of the lack of past data to compare with her

Also, for such an analysis, the number of parameters to take into account is considerable (and tend to infinity)

Can make a parallelism with Chaos Theory and the well-known butterfly effect which have strong applications in weather models

We highly expect the climate change to affect the extreme weather

[extremes in climate change p.347]

It has been proven that winter become warmer in context of RC. (see naveau,...)

?

"The first myth about climate extremes, which has been purported by researchers in climatology or hydrology, among them prominent names, is that "extremes are defined as rare events" or similar. This myth is debunked by a simple bimodal PDF (Fig. 6.12a). The events sitting in the tails of that distribution are not rare" ([Mudelsee, 2014](#), pp.257)

Until now, studies on climate extremes that consider Europe have usually had a strong national signature , or have had to make use of either a dataset with daily series from a very sparse network of meteorological stations (e.g. eight stations in Moberg et al. (2000)) or standardized data analysis performed by different researchers in different countries along the lines of agreed methodologies (e.g. Brazdil et al., 1996; Heino et al., 1999) [Klein Tank et al. \(2002\)](#)

Extrapolation !!!! See p154 [statistical analysis of extreme book]

Voir effet de l'îlot de chaleur -> urbanisation sur les tempêtes !

-> artificial warming on cities stations which were not(less) urbanized 100 years ago.

[In this thesis, efforts have been made to use power of (hyper)references into the text. While this not (yet ?...) usable in printed versions, the reader may feel more comfortable in a numeric version to

more easily handle the vast amount of sections, equations, references, etc... and the links that are made between them.]

We can summarize the research question of this thesis as the following :

-

In this part, we will make use of some general methods to assess if there is indeed a trend in the maximum temperatures

There are two main approaches in EVT, the block-maxima and the peaks-over-threshold approach (see [section 2](#)) yielding to different extreme value distribution. The former aims at while the latter models the (...)

In [chapter 1](#) we will present the method of block-maxima and derive the Generalized Extreme Value (GEV) distributions. In [chapter 2](#) we will . In [chapter 3](#) In chapter 4 we will ..

Finally, we notice that this thesis will more concentrate on the block-maxima (GEV) methods (see [chapter 1](#)), i.e. to data relating to maxima over a period of 1 year. For example here in the introduction, or for the Bayesian analysis in [chapter 5](#).

Part I

Theoretical Framework : Extreme Value Theory

Your writing is quite 'verbose', that is, using many words. Moreover, you tend to repeat yourself. Please write more compactly, more concisely. A written text is not the same as an oral explanation written down.

CHAPTER 1

METHOD OF BLOCK-MAXIMA

block maxima: no hyphen. block-maxima approach: hyphen.

Contents

1.1 Preliminaries	5
1.2 Extremal Types Theorem	7
1.3 Characterization of the GEV distributions : 3 Types	8
1.4 Applications : Examples of Convergence to GEV	10
1.5 Maximum Domain of Attraction	13
1.5.1 Domain of attraction for the 3 types of GEV	14
1.5.2 Closeness under tail equivalence property	17
1.5.3 Domain of attraction of the GEV	18
1.6 The Concepts of Return Levels and Return Periods	18
1.7 Inference	19
1.7.1 Likelihood-based Methods	19
1.7.2 Other Methods	21
1.8 Model Diagnostics : Goodness-of-Fit	21
1.8.1 Return Level Plot	21

In English, it is customary to use capitals: Section 2, Chapter 5, Theorem 3, Definition 1, Figure 6, Appendix A, ...

In this section, we will present the basics of EVT and we will consider a *block-maxima approach*. After defining some useful concepts in [section 1.1](#), we will

~~This chapter is mostly based on~~

For this chapter, we will be mostly based on [Coles \(2001, chapter 3\)](#), [Beirlant et al. \(2006, chapter 2\)](#) and [Reiss and Thomas \(2007\)](#).

1.1 Preliminaries

Why a finite sequence? In limit theorems, it is convenient to have an infinite sequence.

In the following, we will assume a sequence of n independent and identically distributed (iid) random variables that we will write, for convenience, in the form of $\{X_i\} = X_1, X_2, \dots, X_n$. Note that the iid assumption will be relaxed in [chapter 3](#).

Statistical Tools

we let ... denote the i -th ascending order statistic,
First of all, we write the i -th order statistics $X_{(i)}$ which denote the statistics ordered by increasing value

$$X_{(1)} \leq X_{(2)} \leq \dots \leq X_{(n)}, \quad (1.1)$$

assuming we have n observations.

One order statistic is of particular interest for our purpose, the *maximum* $X_{(n)}$

$$X_{(n)} := \max_{1 \leq i \leq n} X_i, \quad (1.2)$$

while for the *minimum* $X_{(1)}$, we can define it with respect to the maximum operator

$$X_{(1)} := \min_{1 \leq i \leq n} X_i = -\max_{1 \leq i \leq n} (-X_i). \quad (1.3)$$

This text will **focus on maxima** but it is important to keep in mind that all the analysis made in the following can be extended to minima through this relation (1.3).
that the analysis through relation

Furthermore, we can easily retrieve the distribution of our statistic of interest $X_{(n)}$, by definition Start new sentence.

$$\begin{aligned} \Pr\{X_{(n)} \leq x\} &= \Pr\{X_1 \leq x, \dots, X_n \leq x\} \\ &\stackrel{(\perp)}{=} \Pr\{X_1 \leq x\} \dots \Pr\{X_n \leq x\} \\ &= F^n(x), \end{aligned} \quad (1.4)$$

where the independence (\perp) follows directly from the iid assumption of the sequence $\{X_i\}$.

First Definitions and Theorems : Motivations

Definition 1.1 (Similar distribution functions). We say that two distribution functions G and G^* are **similar** or are of the **same type** if, for constants $a > 0$ and b we have
"similar"? I don't think this is an "official" term.

$$G^*(az + b) = G(z), \quad \forall z. \quad (1.5)$$

△

It means that the distributions only differ in location and scale. In the sequel, the concept of *similar* distributions will be useful to derive the three different families of extreme value distributions which come from other distributions that are of the *same type*.

What do you mean here by "of the same type"?
As in the definition, or something else?

Principles of stability : Amongst all the principles about EVT that will be covered during this text, the EVT will be highly influenced by the principles of *stability*. It states that a model should remain valid and consistent whatever the choices we make on the structure of this model. For example, if we propose a model for the annual maximum temperatures and another for the 5-year maximum temperatures, the two models should be mutually consistent since the 5-year maximum will be the maximum of 5 annual maxima. Similarly, in an Peaks-Over-Threshold setting (that we will present in 2), a model for exceedances over a high threshold should remain valid for exceedances of higher threshold.

Definition 1.2 (Max-stability). From Leadbetter et al. (1983) or Resnick (1987), we say that a distribution G is **max-stable** if, for each $n \in \mathbb{N}$,

$$G^n(a_n z + b_n) = G(z), \quad n = 1, 2, \dots, \quad (1.6)$$

Don't start a new paragraph. In LaTeX: no empty line in the source code.
for appropriate (normalizing) constants $a_n > 0$ and b_n . \triangle

In other words, taking powers of G results only in a change of location and scale. This concept will be closely connected with the fundamental limit law for extreme values that we will present in the next section. However, the power of max-stable processes is more used in a multivariate setting, whereas we will focus on univariate sequences. Refer for example to Ribatet et al. (2015) for an excellent introduction on max-stable processes. *Min-stability* can easily be found by complement, see for instance Reiss and Thomas (2007, pp.23) for an example.

As this will be fundamental in the following to explain the origin of EVT, we think it is useful to define precisely the concept of *degenerate distribution functions*.

Definition 1.3 (Degenerate distribution functions). We say that the distribution function of a random variable is **degenerate** if it assigns all of the probability to a single point. \triangle

We illustrate this by the construction of the most commonly used theorem in statistics, the Central Limit Theorem (CLT) that we will define below and which typically plays with the empirical mean statistic $\bar{X}_n = n^{-1} \sum_{i=1}^n X_i$. We know that (\bar{X}_n) converges to the true mean μ in probability(?) and thus in distribution, that is to a non-random single point, i.e. to a *degenerate* distribution even almost surely

$$\Pr\{\bar{X}_n \leq x\} = \begin{cases} 0, & x < \mu; \\ 1, & x \geq \mu. \end{cases}$$

That is not very useful, in particular for inferential purposes.

For this reason, CLT aims at finding a non-degenerate limiting distribution for \bar{X}_n , after allowing for normalization by sequences of constants. We will state it in its most basic form :

Theorem 0 (Central Limit Theorem). Let $\{X_i\}$ be the sequence of n iid random variables with $E(X_i^2) < \infty$. Then, as $n \rightarrow \infty$,

$$\sqrt{n}(\bar{X}_n - \mu) \xrightarrow{d} N(0, \sigma^2),$$

where $\mu = E(X_i)$ and $\sigma^2 = V(X) > 0$.

[Note that d means convergence in distribution and the reader may refer to appendix A. for a useful short review of the most important concepts of convergence for EVT.]

Convergence in distribution
is not specific for EVT.

No "."
Whether the review is
useful is up to the
reader to decide :-)
Omit "useful".

Then, by making a proper choice of some normalizing constants, μ and \sqrt{n} (as location and scale parameters respectively), we find the non-degenerate normal distribution in the limit for the empirical mean \bar{X}_n . provided the variance is nonzero

With the same logic, we find for the distribution of maximum order statistics $X_{(n)}$

$$\lim_{n \rightarrow \infty} \Pr\{X_{(n)} \leq x\} = \lim_{n \rightarrow \infty} \Pr\{X_i \leq x\}^n = \begin{cases} 0, & F(x) < 1; \\ 1, & F(x) = 1, \end{cases} \quad (1.7)$$

which is also a degenerate distribution. Thus, this is exactly what Extreme Value Theory also aims to achieve for the, that is finding a non-degenerate distribution in the limit of $X_{(n)}$ by means of normalization. We will see how it works in details in the following section.

1.2 Extremal Types Theorem

Introduced by Fisher and Tippet (1928), later revised by Gnedenko (1943) and finally streamlined by Haan (1970a), the *extremal types theorem* is very important for its applications in EVT. First, recall the distribution of maxima in (1.4). It states the following :

Theorem 1.1 (Extremal Types Theorem). *If the distribution of partial maxima of an iid sequence of R.V. $\{X_i\}$ with common (unknown) distribution F , say, $X_{(n)}$, properly normalized, converges to a non-degenerate limiting distribution G , i.e.*

I don't like the abbreviation R.V. :-)

The sentence is incomplete: If..., then what? What is the conclusion of the theorem?

$$\lim_{n \rightarrow \infty} \Pr\{a_n^{-1}(X_{(n)} - b_n) \leq z\} = F^n(a_n z + b_n) = G(z), \quad \forall z \in \mathbb{R}, \quad (1.8)$$

and for some normalizing constants $a_n > 0$, $b_n \in \mathbb{R}$. □

This theorem considers an iid random sample, but note that it holds true even if the original scheme being no longer independent (we will present the stationary case in section 3.1). Furthermore, we will see in section 1.5 that it actually means that F is in the **domain of attraction** of G where G is called the *Generalized Extreme Value* (GEV) distribution, defined by

$$G(z) = \exp \left\{ - \left[1 + \xi \left(\frac{z - \mu}{\sigma} \right) \right]_+^{-\xi^{-1}} \right\} := G_{\xi, \mu, \sigma}(z), \quad (1.9)$$

where $-\infty < \mu, \xi < \infty$ and $\sigma > 0$ with (μ, σ, ξ) being the three parameters of the model characterizing location, scale and shape respectively. The notation $y_+ = \max(y, 0)$ denotes in the above that $\{z : 1 + \xi \sigma^{-1}(z - \mu) > 0\}$ to ensure the term in the exponential is negative, and hence the distribution function converges to 1. We will keep this notation in the following so it is important to remind that this yields a vital condition for the GEV. In particular, this will define the endpoints from the three different characterizations of this distribution from the values of the shape parameter. This will be detailed in the next section.

From Coles (2001), we introduce an important theorem in Extreme Value Theory and that has many implications. This theorem states the following :

Theorem 1.2. *For any distribution function F ,*

$$F \text{ is max-stable} \iff F \text{ is GEV.} \quad (1.10)$$

□

Hence, any distribution functions that are *max-stables* (recall [definition 1.2](#)) are also GEV which is defined in [??extthm](#)], and vice-versa. To gain interesting insights of the implications of this theorem, we think it is useful to give proof but only for the " \Leftarrow " as the converse requires too much mathematical backgrounds.

Proof : Of which theorem?

- If $a_n^{-1}(X_{(n)} - b_n)$ has the GEV as limit distribution for large n as defined in [\(1.8\)](#), then

$$\Pr\{a_n^{-1}(X_{(n)} - b_n) \leq z\} \approx G(z).$$

Hence for any integer k , since nk is large, we have

$$\Pr\{a_{nk}^{-1}(X_{(n)k} - b_{nk}) \leq z\} \approx G(z). \quad (1.11)$$

- Since $X_{(n)k}$ is the maximum of k variables having identical distribution as $X_{(n)}$,

$$\Pr\{a_{nk}^{-1}(X_{(n)k} - b_{nk}) \leq z\} = \left[\Pr\{a_n^{-1}(X_{(n)} - b_n) \leq z\} \right]^k, \quad (1.12)$$

giving two expressions for the distribution of M_n , by [\(1.11\)](#) and [\(1.12\)](#) :

$$\Pr\{X_{(n)} \leq z\} \approx G(a_n^{-1}(z - b_n)) \quad \text{and} \quad \Pr\{X_{(n)} \leq z\} \approx G^{1/k}(a_{nk}^{-1}(z - b_{nk})).$$

- It follows that G and $G^{1/k}$ are identical apart from location and scale coefficients. Hence, G is *max-stable* and therefore GEV. This gives proof of the **extremal types theorem**, [1.1](#).

There's a big gap in the proof here: you still need the convergence-of-types theorem. Perhaps better don't call this a proof, because it isn't one. You may present the argument as a heuristic one, an intuitive explanation.

Proof of Theorem 1.1 or 1.2? □

In words, it means that taking ^{powers} power of G results only in a change of location and scale, and hence by recalling the expression of the distribution of a maximum (eq.([1.4](#))), it is possible to find the non-degenerate GEV in the limit for this maximum $X_{(n)}$.

1.3 Characterization of the GEV distributions : 3 Types

The shape parameter $\xi \in \mathbb{R}$ is called the *extreme value index* (EVI) and is at the center of the analysis in EVT. It determines, in some degree of accuracy, the type of the underlying distribution. Hence, from this general definition of the GEV distribution in [\(1.9\)](#), we can directly retrieve the **three principal types of EV distributions** from the value ξ :

$$\boxed{\text{I}} \quad \xi \rightarrow 0 : \quad G_1(z) = \exp\left\{-e^{-\frac{z-\mu}{\sigma}}\right\}, \quad -\infty < z < \infty. \quad (1.13)$$

The parameter μ here is not the same as the one in (1.9). Please use a different symbol.

$$\boxed{\text{II}} \quad \xi > 0 : \quad G_2(z) = \begin{cases} 0, & z \leq \mu; \\ \exp \left\{ - \left(\frac{z-\mu}{\sigma} \right)^{-\xi^{-1}} \right\}, & z > \mu. \end{cases} \quad (1.14)$$

$$\boxed{\text{III}} \quad \xi < 0 : \quad G_3(z) = \begin{cases} \exp \left\{ - \left(- \frac{z-\mu}{\sigma} \right)^{-\xi^{-1}} \right\}, & z < \mu; \\ 1, & z \geq \mu, \end{cases} \quad (1.15)$$

where we can see the dependence of the distribution function with respect to the value of the location parameter μ . One can notice that we explicitly retrieve the general definition in 1.9 for type II and type III, that is when $\xi \neq 0$.

Hence, all these three classes of extreme distributions can be expressed in the same functional form as special cases of this single three-parameter ^{distribution}. That is, when $\xi \rightarrow 0$ we retrieve the **type I** or *Gumbel* family (1.13) while $\xi > 0$ and $\xi < 0$ leads to the **type II** or *Fréchet* family and to the **type III** or *Weibull* family, (1.14) and (1.15) respectively.

Density

The density of the GEV distribution (1.9), $g(z) = \frac{dG(z)}{dz}$ (we can assume absolute continuity) can be expressed in two forms, as depicted in table 1.1. This is not an assumption but a fact.

Table 1.1: The two cases for the *density distribution* of the GEV

$\xi \neq 0$	$g(z) = \sigma^{-1} \left[1 + \xi \left(\frac{z-\mu}{\sigma} \right) \right]_+^{-\frac{1}{\xi}-1} \exp \left\{ - \left[1 + \xi \left(\frac{z-\mu}{\sigma} \right) \right]_+^{-\xi^{-1}} \right\};$
$\xi \rightarrow 0$	$g(z) = \sigma^{-1} \exp \left\{ - \left(\frac{z-\mu}{\sigma} \right) \right\} \exp \left\{ - \exp \left[- \left(\frac{z-\mu}{\sigma} \right) \right] \right\}.$

We can now try to visually represent these three families. The following figure 1.1 depicts the GEV, defined with respect to the value of the shape parameter ξ .

We think it is important to point out that the location parameter μ does not represent the mean as in the classic statistical view but does represent the “center” of the distribution, and the scale parameter σ is not the standard deviation but does govern the “size” of the deviations around μ . This can be visualized on the figure C.1 in appendix C where we show the variation of the GEV distribution when we vary these parameters. Also, a Shiny application has been build through the R package to visualize in the best way the influence of the parameters on this distribution (see beginning of chapter 6 for an explanation on the use). We clearly notice that the location parameter only implies a horizontal shift of the distribution, without changing its shape, while we clearly see the influence of the scale parameter ^{spread or dispersion} on the ‘span’ of the distribution around μ . For example if $\sigma \nearrow$, then the density will appear more flat.

In the following, let’s define the *left* and the *right endpoint* of a particular df F , respectively ${}_*x$ and x_* , by :

$${}_*x = \inf\{x : F(x) > 0\}, \quad \text{and} \quad x_* = \sup\{x : F(x) < 1\}. \quad (1.16)$$

A fundamental remark that we must notice is that the Gumbel distribution is unbounded. The Fréchet distribution has a finite left endpoint in ${}_*x = \mu - \sigma \cdot \xi^{-1}$ (the blue circle in figure 1.1), and its upper

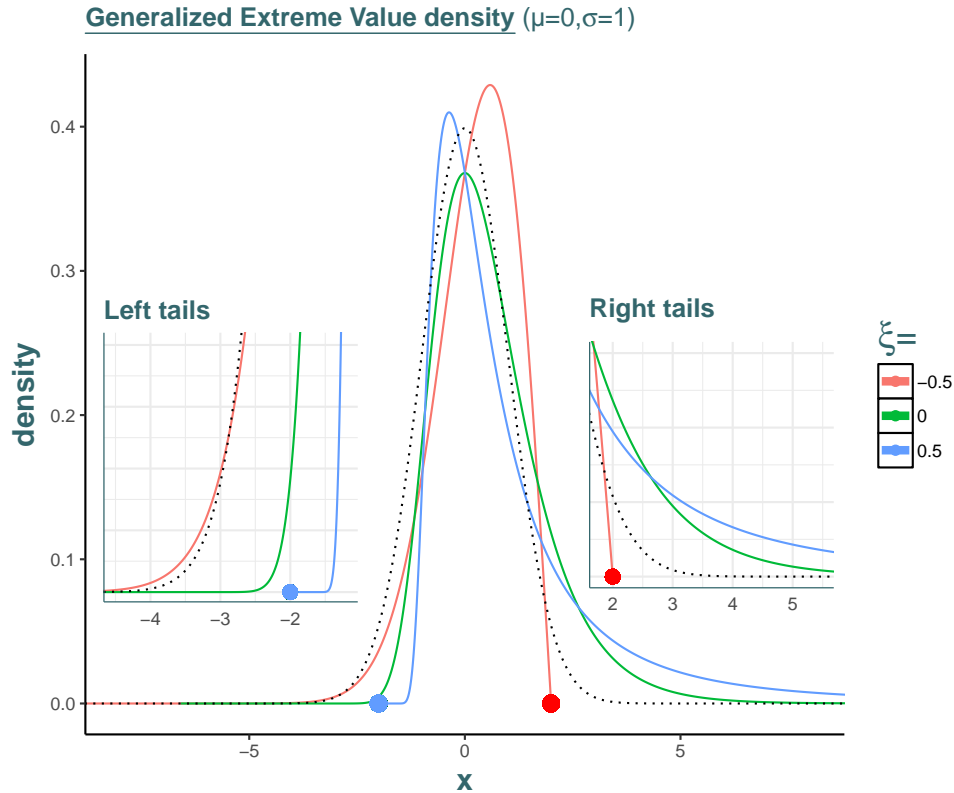


Figure 1.1: GEV distribution with the normal as benchmark (dotted lines) and a zoom on the parts of interest to better visualize the behaviour in the tails. In green, we retrieve the Gumbel distribution ($\xi = 0$). In red, we retrieve the Weibull-type ($\xi < 0$) while in blue, we get the Fréchet-type ($\xi > 0$). The endpoints for the Weibull and the Fréchet are denoted by the red and blue filled circles respectively.

is endpoint tends to $+\infty$ while the Weibull distribution has a finite right endpoint in $x_* = \mu - \sigma \cdot \xi^{-1}$ (the red circle in figure 1.1) and is unbounded in the left. Obviously, this has a serious impact on modelling. This has obviously serious impacts on the modelling. For example, when dealing with maximum temperatures it is intuitive to consider it is more probably bounded to the right while than to left, meaning that a maximum temperature of 70°C is impossible, and hence a Weibull distribution will appear more frequently.

Sometimes, it is useful for people to think not only about the specific form of their data and the distribution they will fit and its characteristics, but also about how we can retrieve these specific distributions in practice. That is the reason why we think it can be useful to detail some examples of how we can construct such extreme distributions for the three types in concrete cases, playing with the appropriate choice of sequences a_n and b_n to retrieve the pertaining distribution family.

Actually, physics says there is an absolute minimum temperature... The presence or absence of finite endpoints is usually not a very important modelling consideration, since things we can measure are bounded anyway.

1.4 Applications : Examples of Convergence to GEV

it is not easy to

In real applications, it is well not easy to find the appropriate sequences. But it is a good exercise to profoundly understand the concept of convergence to GEV by looking at some theoretical examples.

Convergence to Gumbel distribution

The **Type I** or **Gumbel** distribution G_1 can be retrieved by considering, for example, an iid distributed sequence $\{X_j\}$ of n random variables, that is $X_j \stackrel{iid}{\sim} \text{Exp}(\lambda)$ and taking the largest of these values, $X_{(n)}$, as defined earlier. By definition, if the X_j have distribution F , then $F(x) = 1 - e^{-\lambda x}$. Hence, our goal is to find non-random sequences $\{b_n\}$, $\{a_n > 0\}$ such that for positive x

$$\lim_{n \rightarrow \infty} \Pr\{a_n^{-1}(X_{(n)} - b_n) \leq z\} = G_1(z). \quad (1.17)$$

Hence, we can easily find that

$$\begin{aligned} \Pr\{a_n^{-1}(X_{(n)} - b_n) \leq z\} &= \Pr\{X_{(n)} \leq b_n + a_n z\} \\ &= \left[\Pr\{X_1 \leq b_n + a_n z\}\right]^n \\ &= \left[1 - \exp\{-\lambda(b_n + a_n z)\}\right]^n, \end{aligned}$$

from the iid assumption of the random variables and their exponential distribution. Hence, by choosing the sequences $a_n = \lambda^{-1} \log n$ and $b_n = \lambda^{-1}$ and reminding that

you switched the choices for a_n and b_n

$$\begin{aligned} \left[1 - \exp\{-\lambda(b_n + a_n z)\}\right]^n &= \left[1 - \frac{1}{n} e^{-z}\right]^n \\ &\xrightarrow{n \rightarrow \infty} \exp(-e^{-z}) := G_1(z). \end{aligned}$$

Don't overdo the explanations :-)

Recall: $\lim_{n \rightarrow \infty} \left(1 + \frac{x}{n}\right)^n = \exp(x)$

The sentence is incomplete.

Hence, we found the so-called standard *Gumbel* distribution in the limit. Note that the same can be retrieved with $X_j \stackrel{iid}{\sim} N(0, 1)$ and with sequences $a_n = -\Phi^{-1}(1/n)$ and $b_n = 1/a_n$.

Typically unbounded distributions, for example the Exponential and Normal, whose tails fall off exponentially or faster, will have the Gumbel limiting distribution for the maxima. They will have, in particular, medians (and other quantiles) that grow as $n \rightarrow \infty$ at the rate of 'some power of' $\log n$. This is a typical example of light-tailed distribution (i.e., whose tails decay exponentially, as defined in [appendix A.1](#)).

Convergence to Fréchet distribution

The **Type II** or **Fréchet type** (or *Fréchet-Pareto*) distribution $G_2(x)$ has strong relations with the Pareto distribution and also the Generalized Pareto Distribution that will be presented in [chapter 2](#). These are distributions which are typically heavy- or fat-tailed (see [appendix A.1](#)).

Following [Beirlant et al. \(1996\)](#), when starting with a sequence $\{X_j\}$ of n iid random variables following a *basic* (or *generalized* with scale parameter set to 1) Pareto distribution with shape parameter $\alpha \in (0, \infty)$, $X_j \sim Pa(\alpha)$, we have that

$$F(x) = 1 - x^{-\alpha}, \quad x \in [1, \infty). \quad (1.18)$$

Then, by setting appropriately $b_n = 0$, we can write

$$\begin{aligned} -n\bar{F}(a_n z + b_n) &= -n(a_n z + b_n)^{-\alpha} \\ &= \left[F^{\leftarrow} \left(1 - \frac{1}{n} \right) \right]^{\alpha} (a_n)^{-\alpha} (-z^{-\alpha}), \end{aligned}$$

where we define the quantity $F^{\leftarrow}(t) = \inf\{x \in \mathbb{R} : F(x) \geq t\}$ for $t < 0 < 1$ as the *generalized inverse*¹ of F . Hence, it is easy to see that by setting $a_n = F^{\leftarrow}(1 - n^{-1})$ and keeping $b_n = 0$, we have that

$$\Pr\{a_n^{-1}X_{(n)} \leq z\} \rightarrow \exp(-z^{-\alpha}),$$

showing that for those particular values of the normalizing constants, we retrieve the Fréchet distribution in the limit of a basic Pareto distribution. The fact that b_n is set to zero can be understood intuitively since for heavy-tailed distribution (see) such as the Pareto distribution, a correction for location is not necessary to obtain a non-degenerate limiting distribution, see [Beirlant et al. \(1996, pp.51\)](#). More generally, we can state the more general following theorem :

Theorem 1.3 (Pareto-type distributions). *For the same choice of normalizing constants as above, i.e. $a_n = F^{\leftarrow}(1 - n^{-1})$ and $b_n = 0$ and for any $x \in \mathbb{R}$, if*

$$n[\bar{F}(a_n x)] = \frac{\bar{F}(a_n x)}{\bar{F}(a_n)} \rightarrow x^{-\alpha}, \quad n \rightarrow \infty, \quad (1.19)$$

then we say that " \bar{F} is of Pareto-type" or, more technically, " \bar{F} is regularly varying with index $-\alpha$ ". \square

This theorem is interesting to get an understanding of the shape of the tails of this kind of distributions. We ^{define} let the concepts of *regularly varying functions*, together with *slowly varying functions* be defined in [appendix A.3](#)

To avoid confusion with the usual Weibull distribution, one also writes extreme-value Weibull distribution.

Convergence to Weibull distribution

The type **III** or **Weibull** family (?) of distributions $G_3(x)$ arises, for example, in the limit of n iid uniform random variables $X_j \sim U[L, R]$ where L and $R > L$ are both in \mathbb{R} and denote respectively the Left and the Right endpoint of the domain of definition. We have by definition

$$F(x) = \frac{x - L}{R - L}, \quad x \in [L, R].$$

It is 0 for $x < L$ and 1 for $x > R$. We assume we are in the general case, i.e. $[L, R]$ can be $\neq [0, 1]$. Then, we have for the maximum $X_{(n)}$

First define the sequences a_n and b_n

$$\begin{aligned} \Pr\{a_n^{-1}(X_{(n)} - b_n) \leq z\} &= \Pr\{X_{(n)} \leq b_n + a_n z\} \\ &= \left[1 - \frac{R - b_n - a_n z}{R - L} \right]^n, \quad L \leq b_n + a_n z \leq R \\ &= \left(1 + \frac{z}{n} \right)^n \rightarrow e^z, \quad z \leq 0 \quad \text{and} \quad n > |z|. \end{aligned}$$

When choosing $a_n = R$ and $b_n = (R - L)/n$, we find the unit Reversed Weibull distribution $We(1, 1)$ in the limit as expected, that is the Weibull-type GEV with $\xi = -1$. This is a typical example of the

¹from which we can retrieve $x_t = F^{\leftarrow}(t)$, the t -quantile of F . Even if we deal in this text only with continuous and strictly increasing df, we prefer consider *generalized inverse*, for sake of generalization.

maximal behavior for bounded random variables with continuous distributions.

Conditions

Intuitively, it stands to reason that the df F needs certain conditions for the limit to be convergent in [theorem 1.1](#). There exists a *continuity condition* at the right endpoint x_* of F which actually rules out many important distributions. For example, it ensures that if F has a jump at its finite x_* (e.g. discrete distributions), then F cannot have a non-degenerate limit distribution as in [\(1.8\)](#). Examples are well documented in [Embrechts et al. \(2011, section 3.1\)](#) for the Poisson, Geometric and Negative Binomial distributions. We cannot find a nondegenerate distribution in the limit for these distributions even after normalization, which clearly decreases the span .

to exist
What do you mean? The scope for applications?

Comments

Let's finish by noting that it is not mandatory to find the normalizing sequences for inferential purposes, as we will see in [section 3.1](#) by relaxing the independence assumption (which is often poor in practice and too restrictive) to the stationary case. From expression [\(1.8\)](#), we know that

$$a_n^{-1}(X_{(n)} - b_n) \xrightarrow{d} G_{\xi, \mu, \sigma}(z), \quad n \rightarrow \infty. \quad (1.20)$$

After some algebra, we will see that this leads to

$$X_{(n)} \xrightarrow{d} G_{\xi, \mu^*, \sigma^*}(z), \quad n \rightarrow \infty, \quad (1.21)$$

Eq (1.21) is mathematical nonsense :-) The μ^* and σ^* depend on n , so they cannot appear as a limit as n tends to infinity.

where we see that the sequences a_n and b_n have been absorbed into the new location and scale parameters μ^* and σ^* (the shape parameter is invariant). Thus, providing the dependence is limited, we can ignore the normalizing constants in practical applications and fit directly the GEV in our set of maxima. The pertaining estimated parameters μ^* and σ^* will implicitly take the normalization into account. But we will see more details on this at the beginning of [chapter 3](#), and now let's characterize more precisely the distributions pertaining to the GEV.

1.5 Maximum Domain of Attraction

The preceding results can be more easily summarized and obtained when considering *maximum domain of attraction* (MDA). The term "*maximum*" is typically used to make the difference with *sum-stable* distributions but as we only study maxima here, there is no possible confusion in our work. We will then only write "*domain of attraction*" in the following for convenience, considering these two names as synonyms.

The index k is redundant, in view of [\(1.9\)](#).

Definition 1.4 (Domain of attraction). We say that a distribution F is in the **domain of attraction** of an extreme value family $(G_k)_{k=1,2,3}$ in [\(1.13\)-\(1.15\)](#), denoted by $F \in D(G_k)$, if there exist $a_n > 0$ and $b_n \in \mathbb{R}$ such that the distribution of $a_n^{-1}(X_{(n)} - b_n)$ converges in distribution to G_k , where $X_{(n)}$ is the maximum of an iid sequence $\{X_i\}$ with distribution F . \triangle

Let ξ_k denote the EVI pertaining to the extreme value distribution G_k . The above definition is well-defined in the sense that $F \in D(G_i)$ and $F \in D(G_j)$ implies $\xi_i = \xi_j$. One should directly notice the

"The definition is well-defined" ?! :-)

This property is a consequence of the convergence-of-types theorem.

relation with the [theorem 1.1](#).

We have all the necessary tools to the pertaining domains of ^{attraction} attractions now. But, before proceeding, we would like to point out that the fact that the characterization of the first domain of attraction (the Gumbel type) is much more complex than the two following (Fréchet and Weibull class) and requires much more technicalities going beyond the scope of this thesis. ^{Although} Despite this class is important in theory, see e.g. [Pinheiro and Ferrari \(2015\)](#), it is less relevant for our purpose of modelling extremes in a practical case. It often requires other generalizations, for instance with additional parameters to surpass the issues of fitting empirical data. In the last subsection, we will present the unified framework, the domain of attraction pertaining to the GEV distributions, which is a kind of summary for the three first domains of attraction presented.

In each of the characterization of the domains of attractions, we will present some of their most useful, necessary (and sometimes sufficient) conditions. We will especially derive their *von Mises conditions*, coming from [Von Mises \(1936\)](#) but revisited in [Falk and Marohn \(1993\)](#). These conditions are very important in practice and sometimes more intuitive because they make use of the *hazard function* of a df F , defined in the following, for sufficiently smooth distributions :

$$r(x) = \frac{f(x)}{\bar{F}(x)} = \frac{f(x)}{1 - F(x)}. \quad (1.22)$$

It involves the density function $f(x) = \frac{dF(x)}{d(x)}$ in the numerator and it can be thought as a measure of "risk". It can be interpreted as the probability of 'failure' in an infinitesimally small time period between x and $x + \delta x$ given that the subject has 'survived' up till time x .

Please use quotes consistently but sparingly.

1.5.1 Domain of attraction for the 3 types of GEV

Domain of attraction for Gumbel distribution (G_1)

We derive here two ways of formulating necessary and sufficient condition for a distribution function F to be in the Gumbel domain of attraction, namely $F \in D(G_1)$.

The statement of the theorem is incomplete: what is the condition?

Theorem 1.4. Following ([Beirlant et al., 2006](#), pp.72), for some auxiliary function $b(\cdot)$, for every $v > 0$, the condition

$$\frac{\bar{F}(x + b(x) \cdot v)}{\bar{F}(x)} \rightarrow e^{-v}, \quad (1.23)$$

must hold as $x \rightarrow x_*$. Then,

$$\frac{b(x + v \cdot b(x))}{b(x)} \rightarrow 1.$$

□

Why more precise? The above condition is necessary and sufficient, it can't be more precise. Please be more precise in your writing :-)

A lot of more precise characterizations and conditions together with proofs can be found, for example in [Haan and Ferreira \(2006](#), pp.20-33) which is based on the pioneering thesis of [Haan \(1970b\)](#).

Let's now present his **von Mises criterion** as in ([Beirlant et al., 2006](#), pp.73):

Theorem 1.5 (von Mises). If the hazard function $r(x)$ (1.22) is ultimately positive in the neighbourhood

of x_* , is differentiable there and satisfies

$$\lim_{x \uparrow x_*} \frac{dr(x)}{dx} = 0, \quad (1.24)$$

end of proof? sentence incomplete □

then $F \in D(G_1)$. [Reminder: $\lim_{t \uparrow y}(\cdot)$ means that t is approaching y from below, i.e. from values smaller than y in a increasing manner, and vice-versa for $\lim_{t \downarrow y}(\cdot)$]

In words, the slope of the hazard function with respect to x is zero at the limit when x approaches the (infinite) right-endpoint. This ensures a condition on the lightness of the tails of F .

Bad style: paragraph title becomes start of first sentence.

Examples of distributions in $D(G_1)$ include distributions having tails which are exponentially decaying (light-tailed, i.e. the exponential, Gamma, Weibull, logistic, . . .) but also distributions which are moderately heavy-tailed such as the lognormal. To see that, consider a Taylor expansion, we have that

$$\bar{G}_1(x) = 1 - \exp(-e^{-x}) \sim e^{-x}, \quad x \rightarrow \infty,$$

where " \sim " refers to the asymptotic equivalence function. Hence, we directly see the exponential decay of the tails for the Gumbel distribution.

Domain of attraction for Fréchet distribution (G_2)

Let's define $\alpha := \xi^{-1} > 0$ as the *index* of the Fréchet distribution in (1.14).

Definition 1.5 (Power law). *If we look at the tail of the distribution G_2 , a Taylor expansion tells us that*

$$\bar{G}_2(x) = 1 - \exp(-x^{-\alpha}) \sim x^{-\alpha}, \quad x \rightarrow \infty, \quad (1.25)$$

which means that \bar{F}^G tends to decrease as a power law. △

Theorem 1.6. We have *We can say that $F \in G_2$ if and only if*

$$\bar{F}(x) = x^{-\alpha} L(x), \quad (1.26)$$

for some slowly varying function L (see [appendix A.3](#) for the definition). □

In this case and with $b_n = 0$,

$$F^n(a_n x) \rightarrow G_2(x), \quad x \in \mathbb{R},$$

with

$$a_n := F^{\leftarrow}\left(1 - \frac{1}{n}\right) = \left(\frac{1}{1 - F}\right)^{\leftarrow}(n),$$

This previous theorem informs us that all distribution functions $F \in D(G_{2,\alpha})$ necessarily have an infinite right endpoint, that is $x_* = \sup\{x : F(x) < 1\} = \infty$. These distributions are all with regularly varying right-tail with index $-\alpha$ (see [appendix A.3](#)), that is $F \in D(G_{2,\alpha}) \iff \bar{F} \in R_{-\alpha}$.

Finally, let's now present the (revisited) **Von Mises condition** for this domain of attraction which states the following in [Falk and Marohn \(1993\)](#):

If

Theorem 1.7 (von Mises). *if F is absolutely continuous with density f and infinite right endpoint $x_* = \infty$, such that*

$$\lim_{x \uparrow \infty} x \cdot r(x) = \alpha > 0,$$

then $F \in D(G_{2,\alpha})$. □

Explanation not very helpful. Perhaps omit?

In words, it means that when x approaches the (infinite) right endpoint of the distribution and is being multiplied by the hazard function, leads to a non-null constant. This can be thought as a (very small) probability mass remaining even when $x \rightarrow \infty$, and hence a condition on the "heaviness" of the tails. We illustrate this with the standard Pareto distribution, that is

$$F(x) = \left(1 - \left(\frac{x_m}{x}\right)^\alpha\right) 1_{x \geq x_m}, \quad \alpha > 0 \text{ and } x_m > 0.$$

Clearly, we can see that by setting $K = x_m^\alpha$, we obtain $\bar{F}(x) = Kx^{-\alpha}$. Therefore, we have that $a_n = (Kn)^{\alpha^{-1}}$ and $b_n = 0$.

Examples of distributions in $D(G_2)$ include distributions that are typically (very) fat-tailed (or heavy-tailed, see [A.1](#)) distributions, such that $E(X_+)^{\delta} = \infty$ for $\delta > \alpha$. This class of distributions is thus appropriate for phenomena with extremely large maxima, think for example of the rainfall process in some tropical zones. Common distributions include Pareto, Cauchy, Burr, . . . An example to get an idea of this is by looking [\(1.25\)](#) showing that G_2 tends to decrease as a *power law*.

Domain of attraction for Weibull distribution (G_3)

We start by recalling an important relation between the Fréchet and the Weibull distributions

$$G_3(-x^{-1}) = G_2, \quad x > 0.$$

We pointed out the certain 'symmetry' that occurs for these two types (e.g. recall figure [1.1](#)). Hence, this will be useful to characterize $D(G_3)$ using what we know about the Fréchet case.

Theorem 1.8. *We say that $F \in G_3$ as in [\(1.15\)](#) with index $\alpha = \xi^{-1} > 0$ if and only if there exists finite right endpoint $x_* < \infty$ such that*

$$\bar{F}(x_* - x^{-1}) = x^{-\alpha} L(x), \tag{1.27}$$

where $L(\cdot)$ is a slowly varying function (recall [appendix A.3](#)). □

Hence for $F \in D(G_{3,\alpha})$, we have

$$a_n = x_* - F^{\leftarrow}(1 - n^{-1}), \quad b_n = x_*,$$

and hence

$$a_n^{-1}(X_{(n)} - b_n) \xrightarrow{d} G_3.$$

we present

no)

Finally, we still present the **Von Mises condition** related to the G_3 domain of attraction.

Theorem 1.9 (von Mises). *For F having positive derivative on some $[x_0, x_*)$, with finite right endpoint $x_* < \infty$, then $F \in D(G_3)$ if*

$$\lim_{x \uparrow x_*} (x_* - x) \cdot r(x) = \alpha > 0, \quad \int_{-\infty}^{x_*} \bar{F}(u) du < \infty, \quad (1.28)$$

The integral is always infinity,
since $1-F(u)$ converges to 1 as u
tends to $-\infty$? □

Similarly to the Fréchet case, we remark that there is still a probability mass from the hazard rate when x approaches its finite right endpoint, characterized by a non-null constant α which defines the left heavy tail and the right endpoint.

Not all of them.

Examples of distributions in $D(G_3)$ include all the df's that are bounded to the right ($x_* < \infty$). Whereas the Fréchet type is often more preferable in an extreme analysis context because it allows for arbitrarily large values, most phenomena are typically bounded, hence we will think at the Weibull for the most attractive and flexible class for modelling extremes. For example, in our case of modelling a process of (yearly) maximum temperatures, it seems to be the perfect candidate.

1.5.2 Closeness under tail equivalence property

An interesting property of all the three types of domain of attraction $D(G_k)_{k=1,2,3}$ we have derived, is that those are *closed under tail-equivalence*. This is useful for characterizing tail's types of the distributions falling in the pertaining domains of attraction. In this sense,

1. For the **Gumbel** domain of attraction, let $F \in D(G_{1,\alpha})$. If H is another distribution function such that, for some $b > 0$,

$$\lim_{x \uparrow x_*} \frac{\bar{F}(x)}{\bar{H}(x)} = e^b, \quad \text{just } b \quad (1.29)$$

then $H \in D(G_{1,\alpha})$. This emphasizes the exponential type of the tails for the distributions H falling in the Gumbel domain of attraction.

2. For the **Fréchet** domain of attraction, let $F \in D(G_{2,\alpha})$. If H is another distribution function such that, for some $c > 0$,

$$\lim_{x \rightarrow \infty} \frac{\bar{F}(x)}{\bar{H}(x)} = c^\alpha, \quad \text{just } c \quad (1.30)$$

then $H \in D(G_{2,\alpha})$.

3. For the **Weibull** domain of attraction, let $F \in D(G_{3,\alpha})$. If H is another distribution function such that, for some $c > 0$,

$$\lim_{x \uparrow x_*} \frac{\bar{F}(x)}{\bar{H}(x)} = c^{-\alpha}, \text{ just } c \quad (1.31)$$

then $H \in D(G_{3,\alpha})$.

This emphasizes the polynomial types for the tails of the distributions falling in the Fréchet or in the Weibull domain of attraction.

1.5.3 Domain of attraction of the GEV

The conditions that have been stated for the three preceding domains of attraction can be restated under this "unified" framework for the GEV distribution defined in (1.9). For a given df F , by letting the sequences b_n , a_n , and the shape parameter such that

$$b_n = F^{\leftarrow}(1 - n^{-1}), \quad a_n = r(b_n) \quad \text{and} \quad \xi = \lim_{n \rightarrow \infty} r'(x),$$

Then, $a_n^{-1}(X_{(n)} - b_n)$ has the GEV as nondegenerate limiting distribution (differencing the two cases, $\xi \neq 0$ or $\xi = 0$, as done in 1.1 for the density). Among many characterizations, we present the following.
different notation $D(GEV)$? Link to previous theorems?

Theorem 1.10. Let $D(GEV)$ denotes the GEV domain of attraction and F be the df of a sequence $\{X_i\}$ iid. If there exist a positive, measurable function $u(\cdot)$, then for $\xi \in \mathbb{R}$, $F \in D(GEV)$ if and only if:

What is the condition? There do exist positive, mesasurable functions...

$$\lim_{v \uparrow x_*} \Pr \left\{ \frac{X - v}{u(v)} > x \mid X > v \right\} := \lim_{v \uparrow x_*} \frac{\bar{F}(v + x \cdot u(v))}{\bar{F}(v)} = \begin{cases} (1 + \xi x)_+^{-\xi^{-1}}, & \xi \neq 0; \\ e^{-x}, & \xi = 0. \end{cases} \quad (1.32)$$

□

We will see in chapter 2 that it actually defines the "Peaks-Over-Threshold" model.

1.6 The Concepts of Return Levels and Return Periods

After having defined the theoretical properties of distributions pertaining to the GEV family precisely, we are now interested in finding a quantity that could significantly improve the interpretability of such models. *Return levels* play a major role in environmental analysis. For such tasks, it is usually more convenient to interpret EV models in terms of insightful return levels rather than individual parameter estimates.

Assuming for this introductory example our time unit reference is in year -as usually assumed in meteorological analysis-, let us consider the *m-year return level* r_m which is defined as the high quantile for which the probability that the annual maximum exceeds this quantile is $(\lambda \cdot m)^{-1}$, where λ is the mean number of events that occur in a year. λ will be obviously equal to 1 here for yearly blocks, in order to facilitate the interpretation. m is called the *return period* and is, to a reasonable degree of accuracy, the expected time between the occurrence of two so-defined high-quantiles. For example, under stationary assumption, if the 100-year return level is 37 deg c for the sequence of annual maximum temperatures, then 37 deg c is the temperature that is expected to be reached in 100 years. More precisely, you can see it such that r_m is exceeded by the annual maximum in any particular year with probability m^{-1} .

Do not start a sentence with a mathematical symbol.

Let $\{X_{(n),y}\}$ denote the iid sequence of n random variables representing the annual maximum for a particular year y . From (1.9), we have

$$F(r_m) = \Pr\{X_{(n),y} \leq r_m\} = 1 - 1/m$$

$$\stackrel{\text{1-exp?}}{\Leftrightarrow} \left[1 + \xi \left(\frac{r_m - \mu}{\sigma} \right) \right]^{-\xi^{-1}} = \frac{1}{m}.$$

Hence, by inverting this relation, and by letting $y_m = -\log(1 - m^{-1})$, we can retrieve the quantile of the GEV, namely the *return level* r_m

$$r_m = \begin{cases} \mu + \sigma \xi^{-1} (y_m^\xi - 1), & \xi \neq 0; \\ \mu + \sigma \log(y_m), & \xi = 0. \end{cases} \quad (1.33)$$

Henceforth, after having estimated the model (that will be the subject of [section 1.7](#)), we can replace the estimated parameters $\hat{\theta} = (\hat{\mu}, \hat{\sigma}, \hat{\xi})$ in (1.33) to obtain an estimate of the m -year return level.

However, we recall that the definition of return period is easily misinterpreted and the given above is thus not universally accepted. ~~evaporate means something else than what you want to say~~ ^{the one given above} To evaporate this issue, it is important to distinguish stationary from non-stationary sequences. We investigate the return periods and return levels more precisely by relaxing the independence assumption (stationary) and under a climate change environment (nonstationary) in [section 3.4](#). Regarding the diagnostic issues of the model, we present the *return level plot* in [section 1.8.1](#).

1.7 Inference

As we already discussed in (1.20)-(1.21), a great advantage for the modelling of GEV is that we actually do not have to find the normalizing sequences to estimate the parameters of the model. Hence, we will present in this section the main (frequentists) methods of inference for the GEV. These are mostly based on the likelihood ([section 1.7.1](#)) but we will also present other well-known methods that are widely used to estimate GEV parameters like the (probability weighted) moment estimator ([section 1.7.2](#)). Finally, note that there exist estimators for the EVI ξ only, but it is more relevant to leave that for [section 2.5.4](#). After all, we will rely on the Bayesian inference in [chapter 5](#).

1.7.1 Likelihood-based Methods

The most usual method we will first consider is the Maximum Likelihood (ML). Whereas it generally does a good job, it is also very easy and intuitive and to understand and (in general) to implement.

Depicted by [Smith \(1985\)](#), a potential difficulty with the use of likelihood methods for the GEV concerns the regularity conditions that are required for the usual asymptotic properties associated with the maximum likelihood estimator to be valid. Such conditions are not satisfied by the GEV model because the endpoints of the GEV distribution are functions of the parameter values, see the paragraph below (1.16). the following items present the special cases, depending on the value of the EVI ξ .

1. $\xi < -1$: MLE's are unlikely to be obtainable.

2. $\xi \in (-1, -0.5)$: MLE's are generally obtainable but their standard asymptotic properties do not hold.
3. $\xi > -0.5$: MLE's are regular, in the sense of having the usual asymptotic properties.

But fortunately, in practice, the problematic cases in the two first situations ($\xi \leq -0.5$) are rarely encountered for most environmental problems. This situation corresponds to distributions in the Weibull family of GEV with very short bounded upper tail, see for example the red density in figure 1.1 or in figure C.1 (or in the Shiny app) where we better see what defines the borders of the problematic case. The 'bell' of the curve becomes very narrow. If we reach the problematic cases, Bayesian inference which do not depend on these regularity conditions may be preferable. The fact that the distribution of the process of yearly maxima is upper bounded leads us to vastly consider this subject in chapter 5.

Other forms of likelihood-based methods have also emerged to remedy this problem of instability for low values of ξ . Close to a bayesian formulation, **penalized ML** method has been proposed by [Coles and Dixon \(1999\)](#) which adds a penalty term to the likelihood function to "force" the shape parameter to be > -1 , values close to -1 being much larger penalized. We will actually use this method who try to circumvent issues of the usual likelihood computation in section 4.1, for nonstationary sequences.

We are now considering a sequence $\{Z_i\}_{i=1}^n$ of independent R.V. sharing each the same GEV distribution. Let denote $\mathbf{z} = (z_1, \dots, z_n)$ the vector of observations. From the densities of the GEV distribution $g_\xi(z)$ defined in table 1.1, we can derive the log-likelihood $\ell = \log [L(\mu, \sigma, \xi; \mathbf{z})]$, for the two different cases $\xi \neq 0$ or $\xi = 0$ respectively:

1.

$$\ell(\mu, \sigma, \xi \neq 0; \mathbf{z}) = -m \log \sigma - (1 + \xi^{-1}) \sum_{i=1}^n \log \left[1 + \xi \left(\frac{z_i - \mu}{\sigma} \right) \right]_+ - \sum_{i=1}^n \left[1 + \xi \left(\frac{z_i - \mu}{\sigma} \right) \right]_+^{-\xi^{-1}}. \quad (1.34)$$

2.

$$\ell(\mu, \sigma, \xi = 0; \mathbf{z}) = -m \log \sigma - \sum_{i=1}^n \left(\frac{z_i - \mu}{\sigma} \right) - \sum_{i=1}^n \exp \left\{ - \left(\frac{z_i - \mu}{\sigma} \right) \right\}, \quad (1.35)$$

Maximization of this pair of equations with respect to the parameter vector $\theta = (\mu, \sigma, \xi)$ leads to the MLE with respect to the entire GEV family. Note that there is no analytical solution and hence, it must be numerically optimized. Example will be provided in the practical application in section 7.1.1.

From standard MLE theory, we know that the estimated parameter vector $\hat{\theta}$ will be approximately multivariate normal. Inference such as confidence intervals can thus be applied, relying on this approximate normality of the MLE. Hence, problems of this method arise when the approximate normality cannot hold. The underlying inferences will not be sustainable. Whereas [Zhou \(2010\)](#) closed the discussion on the theoretical properties of the MLE, another method is usually more preferable for inference, the *profile likelihood*.

Profile Likelihood

In general, the normal approximation to the true sampling distribution of the respective estimator is rather poor. This is why it is useful to consider an other approach related to the usual likelihood

method, the *profile likelihood* which is often more convenient when a single parameter is of interest. Let's denote it θ_j . Now let's consider the parameter vector $\boldsymbol{\theta} = (\theta_j, \boldsymbol{\theta}_{-j}) = (\mu, \sigma, \xi)$ typically for parameter inference in EVT in a stationary context, where $\boldsymbol{\theta}_{-j}$ corresponds to all components of $\boldsymbol{\theta}$ except θ_j . $\boldsymbol{\theta}_{-j}$ can be seen as a vector of nuisance parameters. The profile log-likelihood for θ_j is defined by

$$\ell_p(\theta_j) = \arg \max_{\boldsymbol{\theta}_{-j}} \ell(\theta_j, \boldsymbol{\theta}_{-j}). \quad (1.36)$$

Henceforth for each value of θ_j , the profile log-likelihood is the maximised log-likelihood with respect to $\boldsymbol{\theta}_{-j}$, i.e. with respect to all other components of $\boldsymbol{\theta}$ but not θ_j . Generalization where θ_j is of dimension higher than one (e.g. in a nonstationary context) is possible.

Another interpretation is related to the χ^2 distribution and the equality with the hypothesis testing the Gumbel case. Details can be found in [Beirlant et al. \(2006, pp.138\)](#).

1.7.2 Other Methods

The Probability-Weighted-Moments Estimator

Introduced by [Greenwood et al. \(1979\)](#), the *Probability-Weighted-Moments* (PWM) of a R.V. X with df F , are the quantities

$$M_{p,r,s} = \mathbb{E} \left\{ X^p [F(X)]^r [1 - F(X)]^s \right\}, \quad (1.37)$$

for real p, r and s . From this equation (1.37), we can retrieve the PWM estimator from specific choices of p, r and s .

1.8 Model Diagnostics : Goodness-of-Fit

After having fitted a statistical model to data, it is important to assess its accuracy in order to infer reliable conclusions from this model. Ideally, we aim to check that our model fits well the whole population, that is the whole distribution of maxima, e.g. all the past and future temperature maxima that will arise. As this cannot be achieved in practice, it is common to assess a model with the data that were used to estimate this model. The aim here is to check that the fitted model is acceptable for the available data.

As these concepts are generally well known in the statistical word, we decide to let in [appendix A.4](#) a reminder of **quantile** and **probability plots** applied in the extremal world.

1.8.1 Return Level Plot

In [section 1.6](#) we introduced the concept of return levels and how it can be useful for intuitive interpretations. Now, we will use this quantity as a diagnostic tool for model checking. Approximate confidence intervals for the return levels can be obtained by the delta method which relies on the asymptotic normality of the MLE, and hence produces a symmetric confidence interval.

Standard errors of the estimates As usual, it is important to compute the standard errors to construct confidence intervals, and hence the return level plot. We naturally expect these standard errors to increase with the return period. Indeed, it is less accurate to estimate return levels for 100 years than for 2 years. As r_m is a function of the GEV parameters, we can use the *delta method* to approximate the variance of \hat{r}_m . Specifically,

$$\text{Var}(\hat{r}_m) \approx \nabla r_m' V \nabla r_m,$$

with V the variance-covariance matrix of the estimated parameters $(\hat{\mu}, \hat{\sigma}, \hat{\xi})'$ and

$$\begin{aligned} \nabla r_m' &= \left[\frac{\partial r_m}{\partial \mu}, \frac{\partial r_m}{\partial \sigma}, \frac{\partial r_m}{\partial \xi} \right] \\ &= \left[1, \xi^{-1}(y_m^{-\xi} - 1), \sigma \xi^{-2}(1 - y_m^{-\xi}) - \sigma \xi^{-1} y_m^{-\xi} \log y_m \right], \end{aligned} \quad (1.38)$$

with $y_m = -\log(1 - m^{-1})$ and the gradient being evaluated at the estimates $(\hat{\mu}, \hat{\sigma}, \hat{\xi})$.

But a problem arise for the so-computed standard errors when considering long-range return levels. They can increase so drastically with the return period that the confidence intervals of the *return level plot* can become difficult to work with. To try to get rid of this issue and to allow for we will construct intervals on the basis of the *profile* log-likelihood. Note that this inference rely on the model adequacy and hence, more uncertainty should be given if the model fit is not perfect.

Profiled likelihood Return levels Usual likelihood methods are not the most accurate for inference in EVT. The problem was that confidence intervals computed in the usual method, with standard errors computed by the Delta method in (1.38), relying on the normal approximation, was not reliable for inference on return levels. This is due to severe asymmetries that are often observed in the likelihood surface for return levels, especially for large quantiles (see Bolívar et al. (2010)).

Profile likelihood method is more accurate for confidence intervals, which better capture the skewness generally associated with return level estimates. We are now specifically interested in computing the profile log-likelihood for the estimation of the return level $\theta_j = r_m$. To do that, we present a method which consists of three main steps :

1. To include r_m as a parameter of the model, by 1.33 we can rewrite μ as a function of ξ, σ and r_m :

$$\mu = r_m - \sigma \xi^{-1} \left[\left(-\log\{1 - m^{-1}\} \right)^{-\xi} - 1 \right].$$

By plugging it in the log-likelihood in (1.34)-(1.35), we obtain the new GEV log-likelihood $\ell(\xi, \sigma, r_m)$ as a function of r_m .

2. We maximise this new likelihood $\ell(\xi, \sigma, r_m = r_m^-)$ at some fixed low value of $r_m = r_m^- \leq r_m^+$ with respect to the "nuisance" parameters (ξ, σ) to obtain the profiled log-likelihood

$$\ell_p(r_m = r_m^-) = \arg \max_{(\xi, \sigma)} \ell(r_m = r_m^-, (\xi, \sigma)).$$

We choose arbitrarily large value of the upper range r_m^+ , and conversely for starting point of r_m^- .

3. Repeat the previous step for a range of values of r_m such that $r_m^- \leq r_m \leq r_m^+$ and then choose

r_m which attain the maximum value of $\ell_p(r_m)$.

By doing this little algorithm, we can easily obtain the *profile log-likelihood plot*.

Interpretation Generally plotted against the return period on a logarithmic scale, the return levels has different shapes depending on the value of the shape parameter ξ , namely :

- If $\xi = 0$, then return level plot will be **linear**.
- If $\xi < 0$, then return level plot will be **concave**.
- If $\xi > 0$, then return level plot will be **convex**.

This can be easily understood as we have seen that $\xi < 0$ implies an upper endpoint and heavy left tail while $\xi > 0$ implies the converse. Henceforth, the "increasing rate" of the return level will decrease as the return period increases for $\xi < 0$ as it cannot go too far away beyond the upper endpoint, and the converse holds for $\xi > 0$.

???[Overfitting problem] (Not to include : to check)

?? A problem of these diagnostics could arise when we focus on prediction accuracy and as we mentioned, the fact that the model is fitted from the data. This well-known problem is called *overfitting*. It can be roughly defined by the process of fitting to noise from the dataset rather than the underlying signal (put ref here). Here, it can be easily explained by the following :

- We are looking for a model which fits the data at best, i.e. for points which are the nearest possible of the diagonal line.
- But, the so-constructed model from which we put the diagnostic is fitted from these original data against which we make the comparison.
- Hence, there could be a incentive to fit a model which fits the most perfectly the available data, that is which points on the diagnostic plots is the nearest possible of the diagonal line. The model is then the best to fit the data at hand
- But, this is a catastrophe when we are seeking at making good predictions from the fitted model, that is making a guess on new, unseen, unavailable data. The model has then lost flexibility, it is not regularized and cannot generalize. (unless the feature space, hear the initial data space, has been completely explored (—>infinite data ?))

See the link with the trade-off bias-variance for threshold selection.

PEAKS-OVER-THRESHOLD METHODS

Contents

2.1 Preliminaries: Intuitions	25
2.2 Characterization of the Generalized Pareto Distribution	26
2.2.1 Outline proof of the GPD and justification from GEV	27
2.2.2 Dependence of the scale parameter σ	28
2.2.3 Three different types of GPD and duality with GEV	28
2.2.4 Examples of the GPD as limiting distribution for exceedances	29
2.3 Return Levels	29
2.4 Point Process Approach	30
2.4.1 Non-homogeneous Poisson Process	30
2.5 Inference	30
2.5.1 Likelihood-based Methods	31
2.5.2 Profile Likelihood	31
2.5.3 Other Methods	31
2.5.4 Estimators Based on Extreme Order Statistics for EVI	31
2.5.5 The Probability-Weighted-Moment Estimator	32
2.5.6 Estimators based on Generalized Quantile	33
2.6 Threshold Selection (Methods)	33
2.6.1 Standard Threshold choice for the excess models	33
2.6.2 "Varying" Threshold : Mixture Models	36
2.6.2.1 Nonstationary extremes	36

Seuils meteo : 0C (gel permanent), 25C et 30C pour les Tx 0C (gel) et 20C pour les Tn

—> Use this for thresholds ?

In this chapter we will focus on the other kind of EV models, modelling excess over a threshold. This... In [section 2.1](#), we briefly introduce the concepts to be able to more precisely characterize the distributions in [section 2.2](#). In [section 2.3](#), we introduce the Poisson models,

2.1 Preliminaries: Intuitions

The threshold models relying on the *Peaks-Over-Threshold* (POT) method are useful to propose a better (?) alternative than the blocking method in 2.1. With this new method, we consider a more natural way of determining whether an observation is extreme or not, by focusing only on all observations that are greater than a pre-specified *threshold*. As we saw, estimates of the GEV parameters are sensitive to the size of block chosen to identify extremes (see) while we will investigate that the estimates of the *Generalized Pareto Distribution* (GPD)¹ parameters are more stable in this sense. Henceforth POT avoids the problem that can arise by considering the maximum of blocks only (), but this method also brings its own problems (). Be aware that this method brings lots of problems with the independence condition... And, especially for temperature data, where for example during heat or cold waves...

Let's consider a sequence $\{X_j\}$ of n iid random variables having marginal distribution function F . We are then regarding for observations that exceed a well-chosen (see) threshold u , which must obviously be smaller than the right endpoint $x_* = \sup\{x : F(x) < 1\}$ of F . The aim here is to find a "child" probability distribution function (fig.? -video youtube), say H , from the underlying (parent) distribution F , that will allow us to model the exceedance $Y = X - u$, and with H then expressed as $H(y) = \Pr\{X - u \leq y \mid X > u\}$. Typically, threshold models can therefore be regarded as the conditional survival function of the exceedances Y , knowing that the threshold u is exceeded (Beirlant et al., 2006, pp.147) :

$$\Pr\{Y > y \mid Y > 0\} = \Pr\{X - u > y \mid X > u\} = \frac{\bar{F}(u + y)}{\bar{F}(u)}. \quad (2.1)$$

or in terms of the exceedance distribution function $F^{[u]}(x) = \Pr\{X \leq u + x \mid X > u\}$ (Reiss and Thomas, 2007, pp.12), ? and Rosso (2015) :

$$F^{[u]}(x) = \frac{\Pr\{X - u \leq x, X > u\}}{\Pr\{X > u\}} = \frac{F(x + u) - F(u)}{\bar{F}(u)}, \quad (2.2)$$

making use of the well-known conditional probability law. One can remark that (2.1) is actually the survivor of the exceedance distribution function, that is $\bar{F}^{[u]}$.

These intuitive characterizations we have given above about the modelling of the threshold exceedances in term of probability distribution function can be useful to understand the following.

However, if the parent distribution F were known, we would be able to compute the distribution of the threshold exceedances in (2.1). (Coles, 2001, pp.74) But as for the GEV in the method of block-maxima (section 2.1), the distribution F is not known in practice, as we will see also in (..). Hence, and as usual in statistics², we must again rely on approximations. This time, we will try to approximate (2.2)

¹Notice that, as an abuse of language and for smoother readability, we will use the abbreviations "GPD" to denote both the Distribution and the Distribution *function*

²?

2.2 Characterization of the Generalized Pareto Distribution

Analogously to the *Fisher-Tippett* theorem in section 2.1 which applies for the block maxima, we have now to define a new theorem which applies for values above a predefined threshold. From this result 1.9(?), these two theorems form together the basis of Extreme Value Theory.

Theorem 2.1 (POT-stability). *Reiss and Thomas (2007, pp.25) The max-stability theorem in ?? can be applied and are formulated here by the fact that the GP distribution functions H are the only continuous one such that, for certain choice of constants a_u and b_u ,*

$$F^{[u]}(a_u x + b_u) = F(x).$$

□

This will be useful for modelling the exceedances in the following theorem (?). And for the examples (see ex. p.25)

Theorem 2.2 (Pickands–Balkema–de Haan). *discovered by Balkema and Haan (1974) and Iii (1975) which showed that the distribution of a threshold u of normalized excesses $F^{[u]}(x)(b_u + a_u x)$, as the threshold approaches the right endpoint x_* of F , is the Generalized Pareto Distribution (GPD) $H_{\xi, \sigma_u}(y)$. That is, if X is a random variable for which (1.8) holds, and for the approximating GP distribution function possessing the same left endpoint u as the exceedance distribution function $F^{[u]}$, we have Reiss and Thomas (2007, pp.27):*

$$|F^{[u]}(x) - H_{\xi, \sigma_u}(x)| \longrightarrow 0, \quad u \rightarrow x_*.$$

□

Or, in an other, maybe more intuitive formulation (the same : delete) Coles (2001) :

$$\Pr\{X(-u) \leq y \mid X > u\} \longrightarrow H_{\xi, \sigma_u}(y), \quad u \rightarrow x_*, \quad (2.3)$$

where the GPD is defined as :

$$H_{\xi, \sigma_u}(y) = \begin{cases} 1 - \left(1 + \frac{\xi y}{\sigma_u}\right)_+^{-\xi^{-1}}, & \xi \neq 0; \\ 1 - \exp\left\{-\frac{y}{\sigma_u}\right\}, & \xi = 0. \end{cases} \quad (2.4)$$

We recall again that $y = x - u > 0$, and where the scale parameter is denoted σ_u to emphasize its dependency with the chosen threshold u :

$$\sigma_u = \sigma + \xi(u - \mu), \quad (2.5)$$

where one can also remark that the location parameter μ does not appear anymore in (??) as it does appear in 2.9.

2.2.1 Outline proof of the GPD and justification from GEV

As we did for block-maxima approach in section 2.1.1 (1.2-1.2), we think it is interesting to have a formal and comprehensive, and still not too technical, intuitive view of where are the GPD from. We remind that we aim here at retrieving the GPD $H_{\xi, \sigma_u}(y)$ (2.3-??) from probability distributions as expressed in (2.1-2.2).

Proof :

- We start with X having distribution function F . From the GEV theorem in section 2.1. (see 1.8-1.9), we have for the largest order statistic, for large enough n ,

$$F_{X_{(n)}}(z) = F^n(z) \approx \exp \left\{ - \left[1 + \xi \left(\frac{z - \mu}{\sigma} \right) \right]^{-\xi^{-1}} \right\}, \quad (2.6)$$

with $\mu, \sigma > 0$ and ξ the GEV parameters. hence, by simply taking logarithm on both sides, we have

$$n \ln F(z) \approx - \left[1 + \xi \left(\frac{z - \mu}{\sigma} \right) \right]^{-\xi^{-1}}. \quad (2.7)$$

- We also have that, from Taylor expansion, $\ln F(z) \approx -[1 - F(z)]$ as both sides go to zero when $z \rightarrow \infty$. Therefore, substituting into (2.7), we get the following for large u :

$$1 - F(u + \mathbf{y}) \approx n^{-1} \left[1 + \xi \left(\frac{u + \mathbf{y} - \mu}{\sigma} \right) \right]^{-\xi^{-1}}.$$

where we specially added the therm $\mathbf{y} > 0$ for our purpose of retrieving something in the form of (2.1)-(2.2).

- Finally, we get for (2.1), with some mathematical manipulations, as $u \rightarrow x_*$:

$$\begin{aligned} \Pr\{X > u + y \mid X > u\} &= \frac{\bar{F}(u + y)}{\bar{F}(u)} \approx \frac{n^{-1} [1 + \xi \sigma^{-1}(u + y - \mu)]^{-\xi^{-1}}}{n^{-1} [1 + \xi \sigma^{-1}(u - \mu)]^{-\xi^{-1}}} \\ &= \left[1 + \frac{\xi \sigma^{-1}(u + y - \mu)}{1 + \xi \sigma^{-1}(u - \mu)} \right]^{-\xi^{-1}} \\ &= \left[1 + \frac{\xi y}{\sigma_u} \right]^{-\xi^{-1}}, \end{aligned}$$

where σ_u is still linear in the threshold u (you will see in (2.9), that is $\sigma_u = \sigma + \xi(u - \mu)$). By simply reverting the probability as in (2.2), we have then

$$\begin{aligned} \Pr\{X - u \leq y \mid X > u\} &= 1 - \Pr\{X > u + y \mid X > u\} \\ &= 1 - \left(1 + \frac{\xi y}{\sigma_u} \right)^{-\xi^{-1}}, \end{aligned} \quad (2.8)$$

which is $GPD(\xi, \sigma_u)$ as required and σ_u is as defined in (2.9).

□

More comprehension can come from (Reiss and Thomas, 2007, pp.27-28) or if one wants to analyse rates of convergence.

2.2.2 Dependence of the scale parameter σ

We chose to express the scale parameter as σ_u to emphasize its dependency with the threshold u . If we increase the threshold, say to $u' > u$, then the scale parameter will be adjusted following :

$$\sigma_{u'} = \sigma_u + \xi(u' - u), \quad (2.9)$$

and in particular, this adjusted parameter $\sigma_{u'}$ will increase if $\xi > 0$ and decrease if $\xi < 0$. If $\xi = 0$, there would be no change in the scale parameter³. We think important to point out the fact that, similarly as mentioned for the GEV models in (1.9), the scale parameter σ_u for GPD models is not the usual standard deviation, but does govern the “size” of the excesses. (AghaKouchak et al., 2013, pp.20)

We will later discuss the threshold choice in section 3.

2.2.3 Three different types of GPD and duality with GEV

One will remark the similarity with the GEV distributions as the parameters of the GPD of the threshold excesses are uniquely determined by the corresponding GEV distribution parameters of block-maxima (see outline proof in the above to convince yourself). Hence, the shape parameter ξ of the GPD is equal to that of the corresponding GEV and, most of all, it is invariant⁴ while the computation of σ_u will not be affected by changes of the corresponding μ or σ in the GEV, from the self-compensation arising in (2.9). (Coles, 2001, pp.76)

Hence, as for the block-maxima approach, there are also three possible families of the GPD depending on the value of the shape parameter ξ which determines the qualitative behaviour of the GPD. Hosking and Wallis (1987), Singh and Guo (1995)

- The **first** type, call it $H_{0,\sigma_u}(y)$, comes by letting the shape parameter $\xi \rightarrow 0$ in ??, giving :

$$H_{0,\sigma}(y) = 1 - \exp\left(-\frac{y}{\sigma}\right), \quad y > 0. \quad (2.10)$$

One can easily notice that it corresponds to an **exponential** distribution function, and hence light-tailed, with parameter $1/\sigma_u$, namely $Y \sim \exp(\sigma_u^{-1})$.

- The **second** and the **third** types, that is when $\xi < 0$ and $\xi > 0$ (resp.), differ only by their support :

³This is consistent with the *memoryless property* of the exponential distribution H_{0,σ_u} (??), for which we give more details in

⁴For instance, choosing different block size in the GEV modelling would shift its (estimated) parameters while GPD (estimated) parameters are *stable*.

$$H_{\xi, \sigma_u}(y) = 1 - \left(1 + \frac{\xi y}{\sigma_u}\right)^{-\xi^{-1}}, \quad \text{for } \begin{cases} y > 0, & \xi > 0; \\ 0 < y < \sigma_u \cdot |\xi|^{-1}, & \xi < 0. \end{cases} \quad (2.11)$$

Therefore, if $\xi > 0$ the corresponding GPD is of **Pareto**-type, hence is heavy-tailed, and has no upper limit while if $\xi < 0$, the associated GPD has an upper bound $y_* = u + \sigma_u/|\xi|$ and is then **Beta**-type distribution. A special case arise when $\xi = -1$ where the pertaining distribution becomes Uniform $(0, \sigma)$. (Grimshaw, 1993, pp.186)

Some plots ?

After looking at the behaviour of the density of these functions, we will procure a more comprehensive view by defining some examples of how to retrieve these different types of Generalized Pareto Distributions.

Density functions of the GPD

$$h_{\xi, \sigma_u}(y) = \frac{\xi}{\sigma_u} \left(1 + \xi \frac{y}{\sigma_u}\right)^{-\xi^{-1}-1} \quad (2.12)$$

2.2.4 Examples of the GPD as limiting distribution for exceedances

We have seen in the previous paragraph that if we can have an approximate distribution G for block-maxima, then threshold excess will have a corresponding distribution given by a member of the Generalized Pareto family. Whence the shape parameter ξ , as for GEV distributions, is determinant for controlling the behaviour of the GPD, and thus leads to the three different types in (2.10)-(2.11).

1. The first type

The choice of a threshold will be discussed in section 3.5.1.

From (Beirlant et al., 2006, p.147-),

2.3 Return Levels

In a similar way as for method of block-maxima (see section 1.6). From (2.4), we obtain the quantiles of the GPD simply by setting this equation equal to $1 - 1/m$ and inverting.

However, differently as for *block*-maxima, the quantiles of the GPD cannot be as readily interpreted as return levels because the observations no longer derive from predetermined *blocks* of equal length. Instead, it is now required to estimate the *probability of exceeding the threshold*, ζ_u .

We can now retrieve the return level r_m , i.e. the **value that is exceeded on average once every m observations**. This value is given by

$$r_m = \begin{cases} u + \sigma_u \xi^{-1} [(m\zeta_u)^\xi - 1], & \xi \neq 0; \\ u + \sigma_u \log(m\zeta_u), & \xi = 0. \end{cases} \quad (2.13)$$

provided m is sufficiently large.

Interpretation

Whereas the interpretation of the plot in function shape parameter value is the same as for the block-maxima method (see the end of [section 1.8](#), it is more convenient to replace the value of m by $N \cdot n_y$ in (2.13), where n_y is the number of observations per year, to give return levels on an annual scale. This method allows us to obtain the *N-year return level* which is now commonly defined as the level expected to be exceeded once every N years.

2.4 Point Process Approach

Following mainly [Coles \(2001\)](#), with some further concepts taken from ? ,

As for the two preceding methods, the point process approach aims at modelling some sequences which are initially assumed to be independent (??)

[coles, pp.124] Here, Point Process could be seen as a kind of summary of the two previous methods (respectively in [chapter 1](#) and in rest of [this chapter](#)), leading to nothing new. However, this approach is often preferred :

1. Its interpretation **unifies** the **models** considered so far.
2. Its likelihood enables a more natural formulation of non-stationarity in excess models from the Generalized Pareto model, see [section 2.2](#).

Furthermore, we will see that the parametrization of the point process model is invariant to threshold choice so that this variation would only affect the well-known (already mentioned) bias-variance trade-off in the inference. Interesting if seasonal modelling.

Hence, we recall that if Y is Poisson distributed with parameter λ , then

$$\Pr[Y = k] = e^{-\lambda} \frac{\lambda^k}{k!}, \quad k \in \mathbb{N}. \quad (2.14)$$

2.4.1 Non-homogeneous Poisson Process

$$N(A) \sim \text{Poi}(\Lambda(A)), \quad \Lambda(A) = \int_A \lambda(x) dx. \quad (2.15)$$

"If a process is stationary and satisfies an asymptotic lack of "clustering" condition for values that exceed a high threshold, then its limiting form is non-homogeneous Poisson with intensity measure" Λ , on a set of the form $A = (t_1, t_2) \times (x, \infty)$, given by

$$\Lambda(A) = (t_2 - t_1) \cdot \left[1 + \xi \left(\frac{x - \mu}{\sigma} \right) \right]_+^{-\xi^{-1}} \quad (2.16)$$

2.5 Inference

We present here some methods that are based on likelihood and other kind methods. We will then more on the Bayesian methods in [chapter x](#)

2.5.1 Likelihood-based Methods

Generalized Pareto Distribution

As we have seen, excess-over-threshold models rely on From (2.12), we can write the *log-likelihood* of the GPD :

$$\ell(\mathbf{z}; \xi, \sigma_u) = -n \ln \sigma_u - (1 + \xi^{-1}) \sum_i^n \ln(1 + \xi \sigma_u^{-1} z_i), \quad (1 + \xi \sigma_u^{-1} z_i) > 0. \quad (2.17)$$

2.5.2 Profile Likelihood

2.5.3 Other Methods

Distinct inference for EVI ξ and global inference (see "others")

Beirlant et al. (2006, pp.140)

As we have seen, the two approaches we have encountered, that is *block-maxima* and POT, share commonly the same parameter ξ . Hence, it is not necessary to differentiate between these methods for the sole estimate the shape parameter.

2.5.4 Estimators Based on Extreme Order Statistics for EVI

These estimators allow to estimate the EVI ξ .

Pickands estimator

Firstly introduced by ?, this method can be applied $\forall \xi \in \mathbb{R}$

$$\hat{\xi}_k^P = \frac{1}{\ln 2} \ln \left(\frac{X_{n-\lceil k/4 \rceil+1,n} - X_{n-\lceil k/2 \rceil+1,n}}{X_{n-\lceil k/2 \rceil+1,n} - X_{n-k+1,n}} \right), \quad (2.18)$$

where we used the definition of Beirlant et al. (2006) We recall that $\lceil x \rceil$ denotes the integer (ceil) part of x .

A condition for the consistency of this estimator is that k must be chosen such that $k/n \rightarrow 0$ as $n \rightarrow \infty$. This condition will hold for the rest rest of the estimators based on (...) in the following

A problem with this intuitive estimator is that its asymptotic variance is very large (see e.g. Dekkers and Haan (1989)) and depends highly on the value of k . To improve this, we can quote the estimator of Segers (2001) which is globally more efficient, depending on the value of an extra-"parameter"(?) and a function to choose.

Methods for heavy-tailed distributions ($\xi > 0$)

Typically, EV analysis of temperature data do not show heavy-tailedness §see). For this reason, some tools commonly used for inference on **Pareto-type** distributions are not relevant. Because of their wide

use and application, we will name them for completeness

The Hill estimator ($\xi > 0$)

This is probably the most simple EVI estimator thanks to the intuition behind its construction. There exists plenty of interpretations to construct it (see e.g. [Beirlant et al. \(2006, pp.101-104\)](#)). Unfortunately, it only holds for heavy-tailed distributions ($\xi > 0$).

It is defined as

$$\xi_k^H = k^{-1} \sum_{i=1}^k \ln X_{n-i+1,n} - \ln X_{n-k,n}, \quad k \in \{1, \dots, n-1\}. \quad (2.19)$$

Following [?](#) , this estimator is consistent. Besides that, this estimator has several problems :

- instability with respect to the choice of k .
- Severe bias due to the heavy-tails of the distribution and thus the slowly varying component which influences negatively.
- Inadequacy with shifted data

Problem : see [pp.105]

The Moment estimator

Introduced by [Dekkers et al. \(1989\)](#), this estimator is a direct generalization of the Hill estimator presented in the previous section.

$$\hat{\xi}_k^M = \hat{\xi}_k^H + 1 - \frac{1}{2} \left(1 - \frac{(\hat{\xi}_k^H)^2}{\hat{\xi}_k^{H(2)}} \right)^{-1}, \quad (2.20)$$

where we define

$$\hat{\xi}_k^{H(2)} = k^{-1} \sum_{i=1}^k (\ln X_{n-i+1,n} - \ln X_{n-k,n})^2.$$

This estimator is also consistent but

Estimator based on generalized quantile plot

To overcome the lack of graphical interpretation of the usual moment estimator,

2.5.5 The Probability-Weighted-Moment Estimator

Probability-Weighted-Moment (PWM) ...

different formulations for POT or block maxima. Look also to "other" directory [Ribereau et al. \(2016\)](#)

The L -Moment Estimator

Wang (1997)

Hosking and Wallis (1997) emphasized the fact that L -moment method came historically as a modification of the PWM method.

2.5.6 Estimators based on Generalized Quantile

2.6 Threshold Selection (Methods)

2.6.1 Standard Threshold choice for the excess models

Single threshold selection involves a **bias-variance trade-off**. That is, (raccourcir)

- **Lower threshold** will induce **higher bias** due to model misspecification. In other words, the threshold must be sufficiently high to ensure that the asymptotics underlying the GPD approximation are reliable.
- **Higher threshold** will induce higher estimation uncertainty, i.e. **higher variance** of the parameter estimate as the sample size is reduced for high threshold.

(Following Leadbetter et al. (1983), this is practically equivalent to estimation of the k^{th} upper order statistic $X_{(n-k+1)}$ called the "tail fraction" below. To ensure tail convergence, as $n \rightarrow \infty$, $k \rightarrow \infty$ but at a reduced rate such that $k/n \rightarrow 0$, i.e. the quantile level of the threshold increases at a faster rate as the sample size n grows.

)

Based on Mean Residual Life

function or *mean excess function*, following again (Beirlant et al., 2006, pp.14-19), (Coles, 2001, pp.78-80),

$$mrl(u_0) := E(X - u_0 \mid X > u_0) = \frac{\int_{u_0}^{x_*} \bar{F}(u) du}{\bar{F}(u_0)}, \quad (2.21)$$

for X having survival function $\bar{F}(u_0)$ computed at u_0 , with $x_* = \sup\{x : F(x) < 1\}$ denoting the right endpoint of the support of F . It denotes, in an actuarial context, the expected remaining quantity or amount to be paid out when a level u_0 has been chosen. However, even if it is mainly applied in an actuarial context or in survival analysis in the literature (see ? for a well-known example), there are also interesting and reliable applications in our more environmental purposes as we will see in the following. Moreover, this function has interesting properties about the tail of the underlying distribution of X (Beirlant et al., 2006, pp.16). In fact, we expect the following :

- If $mrl(u_0)$ is constant, then X has exponential distribution.
- If $mrl(u_0)$ ultimately increases, then X has a heavier tail than the exponential distribution.
- If $mrl(u_0)$ ultimately decreases, then X has a lighter tail than the exponential one.

(and vice-versa, goes it in the two sens??)

This can be particularly interesting for our purpose when considering threshold models. For this case, we can suppose the excesses of a threshold generated by the sequence $\{X_i\}$ follow a generalized Pareto distribution (see 2.2). Knowing the theoretical mean of this distribution, we retrieve, provided the shape parameter $\xi < 1$ and denoting σ_u the scale parameter corresponding to excess of a threshold $u > u_0$,

$$\begin{aligned} mrl(u) &:= E(X - u \mid X > u) = \frac{\sigma_u}{1 - \xi} \\ &= \frac{\sigma_{u_0} + \xi u}{1 - \xi}, \end{aligned} \tag{2.22}$$

from the threshold u dependence with the scale parameter σ (see 2.9). Hence, we remark that $mrl(u)$ is linearly increasing in u , with gradient $\xi(1 - \xi)^{-1}$ and intercept $\sigma_{u_0}(1 - \xi)^{-1}$. Furthermore, we can estimate empirically this function intuitively by

$$\widehat{mrl}(u) = \frac{1}{n_u} \sum_{i=1}^{n_u} (x_{[i]} - u), \tag{2.23}$$

where we let the $x_{[i]}$ denoting the (i-th out of the) n_u observations that exceed u .

Mean residual life plot This leads to an interesting tool for our purpose, the *mean residual life plot*. It comes from combining the linearity detected between $mrl(u)$ and u in (2.22) with (2.23). Therefore, a reliable information can be retrieved from the point of the points

$$\left\{ \left(u, \frac{1}{n_u} \sum_{i=1}^{n_u} (x_{[i]} - u) \right) : u < x_{max} \right\}. \tag{2.24}$$

Even if its interpretation is not easy, this graphical procedure will give insights for the choice of a suitable threshold u_0 to model extremes via general Pareto distribution, that is the threshold u_0 above which we can detect linearity in the plot. Relying on this well-chosen threshold u_0 , the generalized Pareto distribution should be a good approximation. Remind however that its interpretation is often subjective. Furthermore, information in the far right-hand-side of this plot is unreliable. Variability is high due to the limited amount of data (exceedances) above very high thresholds. This can be seen for example on larger confidence intervals.

From (Coles, 2001, pp.83-84)

"Substantial subjectivity in interpreting these diagnostic plots, and the resulting uncertainty. Similar challenges are seen with the River Nidd data, shown in Tancredi et al. (2006), and many other examples in the literature. These examples suggests that a more 'objective' threshold estimation approach is needed and that uncertainty must be accounted for."

see mixture pdf]

Based on the stability of the parameter's estimates

see section 4.3.4 of coles.

montrer pr varying threshold si les estimateurs changent bcp ? Stability plots avec IC grisé qui sajuste complement ds cette region.

The aim is to plot MLE's of the parameters gainst the threshold. These MLE's are supposed to be independent of the threshold choice.

From its simplicity, it forms one of the main tools for the practitioners (as said by e.g.).

But this method is also highly criticized, especially for its lack of interpretability, and the pointwise confidence intervals which are strongly dependent across the range of thresholds (here e.g. we took only....).

Other techniques have thus been proposed, see e.g. [Wadsworth \(2016\)](#) which propose complementary plots with greater interpretability, with a "simple" likelihood-based procedure allowing for automated (more formal ?) threshold selection.

In two words, this method "To identify a treshold that provides the best fit to the likelihood (8)", we maximize the profile likelihood $L_p(j) = L(\hat{\beta}_j, \hat{\gamma}_j, j)$, with $(\hat{\beta}_j, \hat{\gamma}_j)$ the MLE's for a fixed j. After computing $j^* = \operatorname{argmax}_j L_p(j)$, the question is whether $L(\hat{\beta}_{j^*}, \hat{\gamma}_{j^*}, j^*)$ provides a significantly better fit to ξ^* than $L(0, 1, 0) = \prod_{i=1}^{k-1} \phi(\xi_i^*; 0, 1)$. This can be answered by a likelihood ratio test, with test statistic

$$T = \frac{L(\hat{\beta}_{j^*}, \hat{\gamma}_{j^*}, j^*)}{L(0, 1, 0)}. \quad (2.25)$$

If this is significant at level α , there is evidence against a hypothesis of white noise and then we select the threshold $u^* = u_{j^*+1}$ which provides the best fit.

"The lowest threshold that one entertains, u_1 , may also have an impact upon the selected threshold, and might thus be regarded as a tuning parameter. "

"how many thresholds k one should choose. There should be some link to the sample size of the data: if k is too large compared to the sample size n, then the asymptotic theory will not provide a good approximation to the distribution."

Based on the Dispersion Index Plot

As we have seen, the methods considered above lead to a huge amount of subjectivity. Following [Ribatet \(2006\)](#), this method is particularly useful for time series. In [section 2.4](#) we have proven that occurrences of the excesses are represented by a Poisson process, see [2.14](#). Hence, $\mathbb{E}[X] = \operatorname{Var}[X]$ and the *Dispersion Index* statistic introduced by ? is defined by $DI = s^2 \cdot \lambda - 1$, where s^2 is the intensity of the Poisson process and λ can be interpreted as the mean number of events in a block.

A confidence interval can also be computed :

Based on *L-Moments* plot

They are linear combinations of the ordered data values. From the GPD, we have

$$\tau_4 = \tau_3 \cdot \frac{1 + 5\tau_3}{5 + \tau_3}, \quad (2.26)$$

where τ_4 is the *L-Kurtosis* and τ_3 is the *L-Skewness*. See e.g. [Hosking and Wallis \(1997\)](#) for more details on L-moments or [Peel et al. \(2001\)](#) for a known application of this method in hydrology.

We can then construct the *L-Moment plot* :

$$\left\{ (\hat{\tau}_{3,u}, \hat{\tau}_{4,u}) : u \leq x_{\max} \right\} \quad (2.27)$$

where $\hat{\tau}_{3,u}$ and $\hat{\tau}_{4,u}$ are estimations of L-kurtosis and L-skewness based on u and x_{\max} is the maximum observation.

2.6.2 "Varying" Threshold : Mixture Models

[Dey and Yan \(2016\)](#)

see application with pdf small thesis !!! -> unconvulsive.

The threshold is either implicitly or explicitly defined as a parameter to be automatically estimated, and in most cases the uncertainty associated with the threshold choice can be accounted for naturally in the inferences.

The so-called "*fixed threshold approach*" (named in ?, among others) which include thus the diagnostics discussed in [section](#)

There is a wide literature on the subject. The model can be presented in a general way :

$$f(x) = (1 - \phi_u) \cdot b_t(x) + \phi_u \cdot g(x), \quad (2.28)$$

with $b_t(x)$ the density of the bulk model, and where we ignored the parameter dependence for clarity.

"A guiding principle in developing, or choosing, extreme value mixture models is to combine a suitable bulk model, or at least a flexible bulk model, with the tail model. If this is successfully achieved then these models and inference schemes can provide an automated and objective approach to threshold and tail estimation, including uncertainty quantification." book risk pp.62

Problem is the discontinuity which (can) occur in the pdf (not the case for the cdf). this can present bias and uncertainty when the quantity of interest considered is close to the threshold. "Nonstationary extensions of such models can be particularly problematic with the extent of discontinuity varying along the threshold function."

Alternatives are possible to force continuity on the pdf.

"If the bulk model is correctly specified, then the parametric mixture models are easy to understand and quick to fit so are preferred. However, in more usual situation of unknown population distribution, the nonparametric mixture models perform consistently well for low and high quantiles." evmix package thesis.

2.6.2.1 Nonstationary extremes

see thesis2012 p.155

Cross-validation ? Besides all these methods that are very subjective,...

or see gelman bayesian book pp.169

RELAXING THE INDEPENDENCE ASSUMPTION

Contents

3.1 Stationary Extremes	38
3.1.1 The extremal index	40
Clusters of exceedances	40
New parameters	40
Return levels	41
3.1.2 Tail dependence	41
3.1.3 Modelling : Threshold Models	41
3.1.4 Applications	42
3.2 Non-Stationary Extremes	42
3.2.1 Block-Maxima	43
3.2.2 Diagnostics	43
3.3 Model Comparisons	43
3.3.1 Statistical Tools	43
3.4 Return Levels	43
3.5 Inference	44

In environmental applications, the independence assumption is questionable. It is rarely fulfilled , and never completely. From hydrological process (see [Milly et al. \(2008b\)](#) for stationarity) to temperature analysis (see ref) or even the broader area of meteorological applications (see ref), theoretical assumptions that have been made for the models are not sustainable. This sounds also obvious in extreme values analysis of temperature data, since we expect the temperature Whereas it was not really problematic for block-maxima, it was much more painful for POT as we have seen both in section 2. However, here we can see in our case that it does not really happens...(verif.. and why??)

See ([Beirlant et al., 2006](#), pp.375)

....

3.1 Stationary Extremes

From now, we considered $X_{(n)} = \max_{1 \leq i \leq n} X_i$ where we assumed X_1, \dots, X_n are independent random variables. For sake of simplicity, we abandon this notation. In the sequel, this will be denoted

by $\tilde{X}_{(n)} = \max_{1 \leq i \leq n} \tilde{X}_i$ where $\tilde{X}_1, \dots, \tilde{X}_n$ will typically denote a sequence of independent random variables, so that the maximum $\tilde{X}_{(n)}$ is composed of independent random variables only.

We are now interested by modelling $X_{(n)} = \max_{1 \leq i \leq n} X_i$ where $\{X_i\}$ will now denote a *stationary* sequence of n random variables sharing the same marginal distribution as the sequence $\{\tilde{X}_i\}$ of independent random variables, F .

Definition 3.1 (Stationarity). *We say that the sequence $\{X_i\}$ of n random variables is **stationary** if*

More generally, for $h \geq 0$ and $n \geq 1$, the distribution of the lagged random vector $(X_{1+h}, \dots, X_{n+h})$ does not depend on h when the sequence is said to be (strongly) stationary. \triangle

Note that we will only focus on weak(?) stationarity.

For now, we denote $F_{i_1, \dots, i_p}(u_1, \dots, u_p) := \Pr\{X_{i_1} \leq u_1, \dots, X_{i_p} \leq u_p\}$ as the joint distribution function of $(X_{i_1}, \dots, X_{i_p})$ for any arbitrary positive integers (i_1, \dots, i_p) .

Definition 3.2 ($D(u_n)$ dependence condition). *From Leadbetter (1974) and following (Beirlant et al., 2006; Coles, 2001, pp.373-374, pp.93) Let $\{u_n\}$ be a sequence of real numbers. The **$D(u_n)$ condition** holds if for any set of integers $i_1 < \dots < i_p$ and $j_1 < \dots < j_q$ such that $j_1 - i_p > \ell$, we have that*

$$|F_{i_1, \dots, i_p, j_1, \dots, j_q}(u_n, \dots, u_n; u_n, \dots, u_n) - F_{i_1, \dots, i_p}(u_n, \dots, u_n)F_{j_1, \dots, j_q}(u_n, \dots, u_n)| \leq \beta_{n, \ell}, \quad (3.1)$$

where $\beta_{n, \ell}$ is nondecreasing and $\lim_{n \rightarrow \infty} \beta_{n, \ell_n} = 0$, for some sequence $\ell_n = o(n)$, as $n \rightarrow \infty$. \triangle

This condition ensures that, when the sets of variables are separated by a relatively short distance, typically $s_n = o(n)$, the long-range dependence between such events is limited, in a sense that is sufficiently close to zero to have no effect on the limit extremal laws.

From this result, we can retrieve the *extreme-value theorem*

Result is remarkable in the sense that, provided a series has limited long-range dependence at extreme levels ($D(u_n)$ condition makes precise), maxima of stationary series follow the same distributional limit laws as those of independent series. [S.Coles 2001 p.94]

For specific sequence of thresholds u_n that increase with n .

Theorem 3.1 (Limit distribution of maxima under $D(u_n)$). *From Leadbetter (1974). Let $\{X_i\}$ be a stationary sequence of n iid random variables with $X_{(n)} = \max(X_1, \dots, X_n)$. If there exists sequences $\{a_n > 0\}$ and $\{b_n\}$ such that $D(u_n)$ condition holds, then*

$$\Pr\{X_{(n)} \leq u_n\} \longrightarrow H(x), \quad n \rightarrow \infty, \quad (3.2)$$

where H is non-degenerate as defined... and $D(u_n)$ is satisfied with $u_n = a_n x + b_n$ for every real x .

\square

[bootstrap and other...]

Theorem 3.2 (Leadbetter 1983). *From (Coles, 2001, pp.) Let $\{X_i^*\}$ be a stationary sequence and let $\{X_i\}$ be a sequence of iid random variables. By defining $X_{(n)}^* = \max\{X_n^*\}$ and $X_{(n)} = \max\{X_n\}$, we have under regularity conditions,*

$$\Pr\{a_n^{-1}(X_{(n)} - b_n) \leq x\} \longrightarrow G(x), \quad n \rightarrow \infty$$

for normalizing sequences $\{a_n > 0\}$ and $\{b_n\}$, where G is non-degenerate, if and only if

$$\Pr\{a_n^{-1}(X_{(n)}^* - b_n) \leq x\} \longrightarrow G^*(x), \quad n \rightarrow \infty.$$

G^* is the limit distribution coming from a stationary process, defined by

$$G^*(x) = G^\theta(x), \quad (3.3)$$

for some constant $\theta \in (0, 1]$ which is called the **extremal index**.

□

3.1.1 The extremal index

The *extremal index* is an important indicator quantifying the extent of extremal dependence, or equivalently the degree at which the assumption of independence is violated. From eq.(3.3), it is clear that if $\theta = 1$, then the process is independent, but the converse does not hold while the case $\theta = 0$ will not be considered as it is too "far" from independence (check with data?) and brings problems, see for example [Beirlant et al. \(2006, pp.379-380\)](#). Moreover, the results of Theorem 4.2. would not hold true.

.. However, the maximum has a tendency to decrease as ([Coles, 2001, pp.96](#))

Formally, it can be defined as the

$$\theta = \lim_{n \rightarrow \infty} \Pr\{\max(X_2, \dots, X_{p_n}) \leq u_n \mid X_1 \geq u_n\}, \quad (3.4)$$

where $p_n = o(n)$ and the sequence u_n is such that $\Pr\{X_{(n)} \leq u_n\}$ converges. [Coles \(2001\)](#)[slides]

Hence, θ can be thought as the probability that an exceedance over a high threshold is the final element in a *cluster* of exceedances.

Clusters of exceedances

From eq.(3.4), we can now state that extremes have the tendency to occur in cluster, whose *mean cluster size* is θ^{-1} at the limit. Equivalently(?), θ^{-1} is the factor with which the mean distance between cluster is increased.

Identifying clusters and declustering as the distribution of a cluster maximum is the same as the marginal distribution of an exceedance. + slide 82-83(?)

However, [pp.178 Coles], information is discarded when one considers *declustering*. And this information could be substantially important in meteorological applications, for instance to determine heat or cold waves.

New parameters

When $\theta > 0$, we have from Theorem 4.2 that G^* is an EV distribution but with different scale and location parameters than G . If we note by (μ^*, σ^*, ξ^*) the parameters pertaining to G^* and those from G kept in the usual way, we have the following relationships when $\xi \neq 0$

$$\mu^* = \mu - \sigma \xi^{-1} (1 - \theta^\xi), \quad \sigma^* = \sigma \theta^\xi. \quad (3.5)$$

In the Gumbel case ($\xi = 0$), we have $\sigma^* = \sigma$ and $\mu^* = \mu + \log \theta$. The fact that $\xi^* = \xi$ is

Return levels

From that (see clusters), one can see that the probability of an exceedance is variable (see coles, pp.103 or slide 82) (...)

$$r_m = u + \sigma \xi^{-1} \left[(m \zeta_u \theta)^\xi - 1 \right] \quad (3.6)$$

It is important to take that into account as ignorance of this "dependence" can lead to overestimation of the return level.

3.1.2 Tail dependence

From (Reiss and Thomas, 2007, section 2.6), (Coles, 2001, section 8.4) or (Beirlant et al., 1996, section 9.4.1, 10.3.4) + see tail dependence function $\text{atdf}()$ in R.

Problem with traditional tools used in standard time series analysis such as (partial-) autocorrelation functions is that heavy-tailed distributions do not have moments, whereas correlation focus on dependence in the center of the distribution and not the tails. Wada et al. (2016, pp.134) Whence it is important to focus on a *tail dependence measure*.

The auto-tail dependence function using $\chi(u)$ and/or $\bar{\chi}(u)$ employs X against itself at different lags.

a possible estimator (this used by $\text{atdf}()$) can come from the sample version

$$\rho_n(u, h) = \frac{1}{n(1-u)} \sum_{i \leq n} \mathbf{1}(\min(x_i, x_{i+h}) > x_{[nu]:n}) \quad (3.7)$$

(compare with beirlant notations!!!!) $x_{[nu]:n}$

3.1.3 Modelling : Threshold Models

Block-Maxima The modelling with the techniques provided by the GEV distributions (see chapter 1) can be used in the similar way as we have seen from (3.3) or in section 3.1 that the shape parameter remains invariant. The difference is that the effective number of maxima $n = n\theta$ will be reduced and hence the convergence will slower.

A still unsolved problem Coles (2001)[pp.98] is related to the approximations in the limit. Indeed, as the effective number of observations is reduced from n to $n\theta$, the approximation is expected to be

poorer, and this "problem" will be exacerbated with increased levels of dependence in the series.

Thresholds models Practically speaking, one might expect a threshold based analysis to result in estimates of return levels with much reduced standard errors as "all" the extremes are included in the analysis, i.e. those who excess a threshold u . The example in section illustrated this

However, the fact that this method deals with "all" the extremes brings also some problems, and especially the issue of *temporal dependence* (see plot of acf or pacf wrt u) which is illustrated by the fact that the extremes have a tendency to *cluster*. Inferences based on the likelihood found in eq.(2.17) which relies on the independence assumption are now invalid.

Several methods can be used :

- **Filtering out** an (approximate) independent sequence of threshold exceedances.
- **Declustering**. We compute the maximum value in each cluster and then we model these clusters maximums as independent GP random variable. In this approach, we remove **temporal dependence** but we do not estimate it.

However, [Fawcett and Walshaw \(2012\)](#) emphasized the fact that use of the information from *all* extremes rather than just from cluster maxima (?) can be pressed into use to estimate return levels, regardless of the how strong the extremal dependence is. Hence, declustering has no interest. (see book risk p.135) This method accounts for dependence in standard error estimates of the parameters.

3.1.4 Applications

3.2 Non-Stationary Extremes

Whereas we have considered and relaxed during the previous section the first "i" of the "iid" assumption made during the whole chapter 2, we will now tackle the last part "id", i.e. the strong (?) assumption that the observations are **identically distributed**.

The stationarity assumption is very poor to hold for climatologic data [Milly et al. \(2008a\)](#). It is also the case for temperature data.

OUR AIM HERE IS TO MODEL THE EVOLUTION OF As we are dealing with time varying sequences, we can

- Positive trend
- Seasonality

The aim of our modelling will more focus on a different parametrization for the mean, thus in allowing the location parameter to vary through time/seasons.

- Variation in time through t accounting for the season : $\mu(t) = \beta_0 + \mathbb{1}_i(t)$ where $i=1,2,3,4$ represent the seasons.

3.2.1 Block-Maxima

As we continue to consider modelling as yearly blocks, we do only face nonstationary concerns for the trend which is (probably) imputed to the Global Warming. The evidence of seasonality arising when we decrease the length of the blocks is not an issue for yearly modelling. However, we loose information, or comparatively, we do not use all the information as at least one half

3.2.2 Diagnostics

Gumbel plot (slide 94) coles

3.3 Model Comparisons

3.3.1 Statistical Tools

In order to compare our models, that is for example to check whether a trend (or seasonality) is statistically significant, or if the nonstationnary models provide an improvement over the simpler (stationary) model, we will make use of the **deviance statistic** defined as

$$D = 2\{\ell_1(\mathcal{M}_1) - \ell_0(\mathcal{M}_0)\}, \quad (3.8)$$

for two nested models $\mathcal{M}_0 \subset \mathcal{M}_1$, where $\ell_1(\mathcal{M}_1)$ and $\ell_0(\mathcal{M}_0)$ are the maximized log-likelihoods under models \mathcal{M}_1 and \mathcal{M}_0 respectively as defined in .

Asymptotically, the distribution of D is χ_k with k (df) representing the difference of parameters between model \mathcal{M}_1 and \mathcal{M}_0 . Thus, comparisons of D with the critical values from χ_k will guide our decision.

3.4 Return Levels

Stationarity Under an assumption of a stationary sequence, the return level is the same for all years, and this gives rise to the notion of the return period (or m -year event). Hence, the return period of a particular event is the inverse of the probability that the event will be exceeded in any given year. The m -year return level is associated with a return period of m years. However, there are two main interpretations in this context for return periods.

(?, pp.100)

Denoting $X_{(n),y}$ the annual maximum for year y . Omitting the notational dependence on block size n , we assume $\{X_{(n),y}\} \stackrel{iid}{\sim} F$.

1. The first interpretation of the m -year event is **the expected waiting time until an exceedance occurs**. To see that, letting T be the year of the first exceedance, we we can write

$$\begin{aligned}
\Pr\{T = t\} &= \Pr\{X_{(n),1} \leq r_m, \dots, X_{(n),t-1} \leq r_m, X_{(n),t} > r_m\} \\
&= \Pr\{X_{(n),1} \leq r_m\} \dots \Pr\{X_{(n),t-1} \leq r_m\} \Pr\{X_{(n),t} > r_m\} && [\text{iid assumption}] \\
&= \Pr\{X_{(n),1} \leq r_m\}^{t-1} \Pr\{X_{(n),1} > r_m\} && [\text{stationarity}] \\
&= F^{t-1}(r_m)(1 - F(r_m)) \\
&= (1 - 1/m)^{t-1}(1/m).
\end{aligned} \tag{3.9}$$

We easily recognize that T has geometric density with parameter $1/m$. From simple properties of geometric distributions, we found its expected values is $1/(1/m)$, that is the expected waiting time for an m -year event is m years.

2. The second interpretation of the m -year event is that **the expected number of events in a period of m years is exactly 1**. To see that, we define

$$N = \sum_{y=1}^m I(X_{(n),y} > r_m),$$

as the random variable representing the number of exceedances in m years (where I is indicator function). We can view each year as a "trial", and from the fact that we have assumed $\{X_{(n),y}\}$ are iid, we can compute the probability that the number of exceedances in m -years is k

$$\Pr\{N = k\} = \binom{m}{k} (1/m)^k (1 - 1/m)^{m-k},$$

from which we recognize a well-know distribution, that is $N \sim \text{Bin}(m, 1/m)$. Again from properties of this distribution, we easily find that N has an expected value of 1.

Non-stationarity From the definition on non-stationary process, the modelling of return period will change over time. Hence, we introduce the notation of the distribution function F_y of a particular $X_{(n),y}$. We must study

$p(y) = \Pr(X_{(n),y} > r) = 1 - F_y(r)$. If we estimate F_y , we can retrieve easily $p(y)$. $F_y(r_p(y)) = 1 - p$ with the exceedance level $r_p(y)$ changing with year. It shows the changing nature of "risk".

Return period as expected waiting time

Return period as expected number of events

3.5 Inference

NEURAL NETWORK (AND OTHERS)

Contents

4.1 Improvements For Modelling Non-stationary Sequences	45
4.1.1 Generalized Likelihood Methods	45
4.2 Neural-Network Based Inference	45
4.3 Bagging	47
4.4 Bootstrap Methods in EVT	47
4.4.1 Moving Block Bootstrap	48
4.5 Markov models	48

In the era of Artificial Intelligence (AI) and Machine Learning, or more precisely of Artificial Neural Networks (ANN) or the trendy term ‘deep learning, it is interesting to study how this complex mechanism works and how it can be efficiently applied to deal with the issue of nonstationarity in EVT that we confront.

4.1 Improvements For Modelling Non-stationary Sequences

4.1.1 Generalized Likelihood Methods

introduced by ? As we have seen in [section 4.1](#), " The ML method may diverge when sample size is small. To resolve the problems of divergence occurring in the numerical techniques used for ML, Martins and Stedinger [2000] suggest the use of a prior distribution for the shape parameter of the GEV model such that the most probable values of the parameter are included"

When dealing with nonstationnary processes, it is interesting to consider Generalized Maximum Likelihood (GML) estimators. In this case, ? have proven that GML is likely to outperform the usual ML inference.

The GML estimator corresponds to the mode of the empirical posterior distribution...

Properties of the GML estimator see pp740. ?

4.2 Neural-Network Based Inference

Parametric or nonparametric ? It is always a difficult task to state whether Neural Networks (NN) are parametric models or not. NN are somewhere in the gray area between a *parametric* and a *non-*

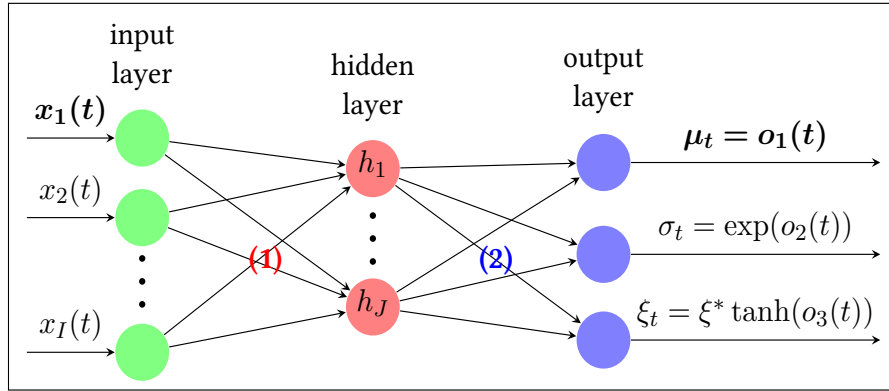


Figure 4.1: Neural Network applied to GEV. Figure made with *tikzpicture* and based on Cannon (2010)

parametric model, in the sense that it assumes the GEV distribution from the output layer which are defined by the three parameters of interest, while it also allows for a fabulous flexibility coming from the hidden layers and which lead to think that these are rather nonparametric. However, this terminological question is actually not relevant here and it is more important to focus on the model of NN that is employed here.

have the power to manage several outputs in a

"Model parameters are estimated via the GML approach using the quasi-Newton BFGS optimization algorithm, and the appropriate GEV-CDN model architecture for each location is selected by fitting increasingly complicated models and choosing the one that minimizes appropriate cost-complexity model selection criteria. For each location examined, different formulations are tested with combinational cases of stationary and nonstationary parameters of the GEV distribution, linear and nonlinear architecture of the CDN and combinations of the input covariates "

? enlightens the following : Provided enough data, hidden units and an appropriate optimization, the NN can capture any smooth dependencies (relationships) of the parameters on the input, i.e., given the input, it can theoretically capture any conditional continuous density, be it asymmetric, **multi-modal**, or heavy-tailed.

(see Cannon (2010) just before conclusion) One could for example expect to have particular relationships between the covariate (time or) and the parameters of interest. Only considering a linear or quadratic trends in the location parameter μ (ore more ? see section. see for other parameters) could thus be seen as a weak modelling procedure, especially when we assumed no reliable prior knowledge on the subject (see bayesian section - hyperref it). NN models have this facility of being capable of modelling any relationships without explicitly specify it *a priori*. To model correctly, one should be able to explicitly discover particular patterns (e.g., which nonlinear or linear relationship between time and TN). This is avoided here because this is done automatically through the NN process.

Physical process such as temperature or even other meteorological data (rainfall as demonstrated by Cannon (2010),...) have this tendence of demonstrating nonlinearities (see ref?) and so are NN's interesting.

As we mentioned, the NN is meant to approximate any functions with good accuracy. It comprise thus all the models considered so far, such as linear tren din μ , quadratic, etc...

Cannon (2010) recommended to use between 1 and 3 (4) hidden layers due to the relatively small

sample of annual extremes (here 117).

From this, we must pay attention to the high danger of *overfitting* (see) which occurs for this sort of models. The other pitfall is its lack of interpretation of the retrieved relationships.

"It bears noting that sensitivity analysis methods, for example, the one used by Cannon and McKendry (2002), are applicable to CDN models and could be used to identify the form of nonlinear relationships between covariates and GEV distribution parameters or quantiles."

... Intro to deep learning.. ?

4.3 Bagging

Nowadays, bagging is used in many state-of-the-art algorithms such as Random Forests (see .. for comparisons of such techniques) and is one of the *ensemble methods* which are praised in Machine Learning for their performance.

In a "pure" climatological point of view, *ensemble models* are of major utility, especially to make weather forecasting, see for example [Suh et al. \(2012\)](#) or, among others

For our purpose, we present another kind of ensemble modelling which is *bagging*.

" Bagging (stands for Bootstrap Aggregation) is the way decrease the variance of your prediction by generating additional data for training from your original dataset using combinations with repetitions to produce multisets of the same cardinality/size as your original data. By increasing the size of your training set you can't improve the model predictive force, but just decrease the variance, narrowly tuning the prediction to expected outcome "

"model averaging, which involves taking a weighted average of multiple models, has been recommended as a means of improving estimation performance Burnham and Anderson, 2004. This approach has been applied successfully in the context of CDN models by Carney et al. (2005) and is worth exploring for GEV-CDN models."

?, pp.256-267 (deep learning html book) The individual classifiers' predictions (having equal weightage) are then combined by taking majority voting. This typically reduces the variance and then the (possible) overfitting

4.4 Bootstrap Methods in EVT

" like the estimated parameters themselves, the SE may not be reliable for small samples. One way to tackle this problem and improve the accuracy of SE is through the bootstrap technique (Efron, 1979). The scheme for using this technique for EV distribution function is described in detail by Katz et al. (2002), including for nonstationary cases, in which the bootstrap samples are manufactured through Monte Carlo resampling of residuals (Equation (8)) to attend to the underlying assumption that original sample consist of iid data. Following this procedure, the bootstrap procedure was designed for generating 1000 samples from each original sample, considering whole year as a bloc"

In [Cannon \(2010\)](#), "he parametric bootstrap outperformed the residual bootstrap" Moreover, " It is possible that alternative bootstrap approaches, for example, the bias-adjusted percentile estimators evaluated by Kysely (2008), might yield better calibrated confidence intervals, although improvements were modest for stationary GEV models. "

For confidence intervals : see [Cannon \(2010, pp.681\)](#) following these steps

1. Fit a nonstationary model to the data
2. Transform the residuals from the fitted model so that they are identically distributed :

$$\varepsilon_t = \left[1 + \xi_t \sigma_t^{-1} (y_t - \mu_t) \right]^{-\xi^{-1}} \quad (4.1)$$

3. etc..

Monte-Carlo based methods, same as Bayesian.

Study and comparisons on the performance (coverage,..) of the methods used for the CI (boot, bayesian, likelihood, asymptotics,...)

4.4.1 Moving Block Bootstrap

[Bootstrap and other resampling in pp.13]

4.5 Markov models

book risk pp.136, [Shaby et al. \(2016\)](#) + code

BAYESIAN METHODS

Contents

5.1 Prior Elicitation	50
5.1.1 Non-informative Priors	50
5.1.2 Informative Priors	51
5.2 Bayesian Computation : Markov Chains	52
5.2.1 Algorithms	52
5.2.2 Hamiltonian Monte Carlo	53
5.2.3 Computational efficiency comparison	53
5.3 Convergence Diagnostics	53
5.3.1 Proposal Distribution	53
5.3.2 The problem of auto and cross-correlations in the chains	54
5.4 Posterior Predictive	54
5.5 Bayesian Predictive Accuracy for Model Validation	55
5.5.1 Cross-validation for predictive accuracy	55
5.6 Bayesian Inference ?	57
5.6.1 Bayesian Credible Intervals	57
5.6.2 Distribution of Quantiles : Return Levels	57
5.7 Bayesian Model Averaging	57

see evdbayes pdf package r

Attention : π ou f ??????

We let useful relevant tools regarding bayesian inference in appendix A.3

Definition 5.1 (Posterior distribution). Let $\mathbf{x} = (x_1, \dots, x_m)$ denote the observed data of a random variable X distributed according to a distribution with density function

$$\pi(\theta|\mathbf{x}) = \frac{\pi(\theta) \cdot L(\theta|\mathbf{x})}{\int_{\Theta} \pi(\theta) \cdot L(\theta|\mathbf{x}) \cdot d\theta} \propto \pi(\theta) \cdot L(\theta|\mathbf{x}) \quad (5.1)$$

where $L(\cdot)$ denotes the likelihood function, as in ?? but there it is the log-likelihood !!! and θ usually denotes the multidimensional set of parameters in EVT, $\theta = (\mu, \sigma, \xi)$, at least in a univariate stationary context. \triangle

1. Whenever it is possible, it allows to introduce to introduce other source of knowledge coming from the domain at-hand, by the elicitation of a prior. The counter-argument of this advantage is that it also introduces (improper ?) subjectiveness.

2. "account- ing for parameter and threshold uncertainty is perhaps handled most easily in the Bayesian paradigm" (Dey and Yan, 2016, pp.106)

As such, It permits an elegant way of making future predictions which is one of the most(?) important issue in EVT.

3. Bayesian framework can overcome the regularity conditions of the likelihood inference (see [section 4.1](#)). Thus it usually provides a viable alternative in cases when MLE (for example) breaks down. And actually, we are not so far from the problematic situations depicted in [section 3.1](#). Moreover, the Highest Posterior Probability (HPD) region is constructed so that... and there is no more need to fall to asymptotic theory as in conventional methods.
4. For an asymmetric distribution, the HPD interval can be a more reasonable summary than the central probability interval (see illustration ...). For symmetric densities, HPD and central intervals are the same while HPD is shorter for asymmetric densities. See [Liu et al. \(2015\)](#)...

As the dependence becomes stronger, the run length n must be larger in order to achieve the same precision. Dependence exists both within the output for a single parameter (autocorrelations) and across parameters (cross-correlations), we discuss this issue in [section text](#).

5.1 Prior Elicitation

Sometimes viewed as advantage from the amount of information that can be retrieved, and sometimes viewed as an drawback due to the (rather unquantifiable) subjectivity that introduced, the construction of the prior is a key step in Bayesian analysis.

Priors are necessary in the Bayesian paradigm to be able to compute the posterior in [\(5.1\)](#). But, priors require the legitimate statement of domain's expert, to make this viewed the less subjective as possible

Prior may not be of great importance if the size (m) of the dataset is large. It can be seen from [\(5.1\)](#) where the amount of information contained in the data through $L(\theta|\mathbf{x})$ will be prominent compared to this contained in the prior through $\pi(\theta)$. Prior will have limited influence.

One is aware that this is not often the case in EVT cases. By design, we are dealing with rather small so-constructed datasets. And mostly for this reason, it could be important to incorporate additional information in this limited dataset through the prior distribution.

5.1.1 Non-informative Priors

Receive a correct and accepted advice from an expert is often difficult So, in many cases, we cannot inject information through the prior. We must then construct a prior which represent this lack of knowledge so that they do not influence posterior inferences.

There exists a vast amount of uninformative priors in the literature (see e.g. [Yang and Berger \(1996\)](#), [Ni and Sun \(2003\)](#)) This family of priors can be *improper*, i.e. priors for which the integral of $\pi(\theta)$ over the parameter space is not finite. It is valid to use improper priors only if the posterior target is proper.

Adjustments of these priors must always be thought in practical applications

Jeffrey's prior

is specified as

$$\pi(\theta) \propto \sqrt{\det I(\theta)}, \quad \text{where} \quad I_{ij}(\theta) = \mathbb{E}_{\theta} \left[- \frac{\partial^2 \log f(\mathbf{x}|\theta)}{\partial \theta_i \partial \theta_j} \right], \quad i, j = 1, \dots, d. \quad (5.2)$$

where $f(\mathbf{x}|\theta)$ is of course the density function of X .

This prior is invariant to reparametrization, but has complex form for EV models, and it exists only when $\xi > -0.5$ in GEV models, where it is function of ξ and σ only.

MDI prior

Maximal Data Information priors

However, it has been showed by [Northrop and Attalides \(2016\)](#) that both Jeffrey and MDI priors give improper posterior when there are no truncation of the shape parameter, i.e. we must restrict the fact that $\pi(\theta) \rightarrow \infty$ as $\xi \rightarrow (-)\infty$ for Jeffreys (MDI), in order to obtain a proper posterior.

Vague priors

The last and often preferred alternative to construct uninformative priors is to use proper priors which are near flat, e.g. which are uniform or with exhibits large variance for the normal distribution.

In GEV we will take independent normal-distributed priors each with a large (tuned) variance. When these variances increase, we get at the limit

$$\pi(\theta) = \pi(\mu, \nu, \xi) \stackrel{(\perp)}{=} \pi(\mu) \cdot \pi(\nu) \cdot \pi(\xi) \propto 1, \quad (5.3)$$

where $\nu = \log \sigma$.

Taking multivariate normal distribution as prior has also been proposed (see) is often difficult as it involves 9 (hyper)parameters in total and this can be difficult

5.1.2 Informative Priors

STAN : "It can also be a huge help with computation to have less diffuse priors, even if they're not informative enough to have a noticeable impact on the posterior. "

Gamma Distributions for Quantile Differences**Beta Distributions for Probability Ratios**

The Bayes Factor

5.2 Bayesian Computation : Markov Chains

Methods have been developed for sampling from arbitrary posterior distributions $\pi(\theta|\mathbf{x})$. Simulations of N values $\theta_1, \theta_2, \dots, \theta_N$ that are iid from $\pi(\theta|\mathbf{x})$ can be used to estimate features of interest.

But simulating from $\pi(\theta|\mathbf{x})$ is usually not achievable and this is why we need **Markov Chain Monte Carlo** (MCMC) techniques. We use it to simulate a markov chain $\theta_1, \theta_2, \dots, \theta_N$ that conerge to the target distribution $\pi(\theta|\mathbf{x})$. This means that, after some *burn-in period* B , $\theta_{B+1}, \dots, \theta_N$ can be treated as random sample from $\pi(\theta|\mathbf{x})$.

Let's now (a bit weakly) define one of the most important results in Markov Chain theory.

Definition 5.2 (*First-order discrete-time Markov Property*). *Let k_0, k_1, \dots be the states associated to a sequence of time-homogeneous random variables, say $\{\theta_t : t \in \mathbb{N}\}$. The Markov property states that the distribution of the future state θ_{t+1} depends only on the distribution of the current state θ_t . In other words, given θ_t , we have that θ_{t+1} is independent of all the states prior to t . We can write this as*

$$\Pr\{\theta_{t+1} = k_{t+1} \mid \theta_t = k_t, \theta_{t-1} = k_{t-1}, \dots\} = \Pr\{\theta_{t+1} = k_{t+1} \mid \theta_t = k_t\}. \quad (5.4)$$

△

or see [Angelino et al. \(2016, section 2.2.3\)](#) for more in-depth results.

The samples are not independent, and the dependence influences the accuracy of the posterior estimates. As dependence becomes stronger, we must increase the run-length N to achieve the same accuracy.

5.2.1 Algorithms

We are looking for a so-generated chain that has a stationary distribution $\pi(\theta|\mathbf{x})$. This is the case if the chain is

1. *aperiodic*
2. *irreducible* or *ergodic*, that is if any state for θ can be reached with probability > 0 in a finite number of steps from any other state for θ .

"The Markov chains Stan and other MCMC samplers generate are *ergodic* in the sense required by the Markov chain central limit theorem, meaning roughly that there is a reasonable chance of reaching one value of theta from another." ?

With MH or Gibbs sampler, we need to tune individually the proposal standard deviations to reach a correct acceptance, and this is often done with trial-and-error methodology.

The performance of the standard Markov chain Monte Carlo estimators depends on how effectively the Markov transition guides the Markov chain along the neighborhoods of high probability. If the exploration is slow then the estimators will become computationally inefficient, and if the exploration is incomplete then the estimators will become biased [Betancourt \(2016\)](#). It is then necessary to consider other form of sampling...

5.2.2 Hamiltonian Monte Carlo

Package Rstan

Neal and others (2011) and Betancourt and Girolami (2015) are really

HMC permit to better exploit the properties of the target distribution to make informed jumps through neighborhoods of high probability while avoiding neighborhoods of low probability entirely.

5.2.3 Computational efficiency comparison

In modern statistical area, computing methods have been widely ... And this need for computations will rise in the future.

We will then compare our 3 methods too see if effectively

5.3 Convergence Diagnostics

When applying MCMC algorithms to estimate posterior distributions, it is vital to assess convergence of the algorithm to try to ensure that we reached the stationary target distribution. Let's now enumerate some of the key steps we must keep in mind when thinking about convergence, and hence reliable results.

1. A sufficient *burn-in period* $B < N$ must be chosen to ensure that the convergence to the posterior distribution $\pi(\theta|\mathbf{x})$ has occurred.
2. For the same reason, a sufficient number of simulations N to eliminate the influence of initial conditions and ensure accuracy in the estimations ((and then make sure than we are sampling from the target stationary (posterior) distribution)).
3. Several dispersed starting values must have been simulated to ensure we explored all the regions of high probability. This is particularly important when the target distribution is complex.
4. The chains must have good mixing properties, in the sense that the whole parameter space (...) A common technique that we will apply is to run different chains several times and then combine a proportion of each chain (typically 50%) to get the final chain. This procedure wants to ensure a proper mixing behaviour. The potential scale reduction factor (Gelman diagnostic) is also a popular tool, see .

We must keep in mind that no convergence diagnostics can prove that convergence really happened and validate the "model". However, a combined use of several relevant diagnostics will be required to increase our confidence that convergence actually happened.

5.3.1 Proposal Distribution

The main ideas are :

- If the variance of the proposal distribution is too large, most proposals will be rejected :
ie the jumps through the chain are too large,

- If the variance of the proposal distribution is too low, then most proposals will be accepted

Both are harmful for the objective of an efficient "visit" of the whole parameter space.

Widely speaking, we consider 2 different types of algorithms in which it is preferable to target a certain acceptance rate. It is distinguished by the updating manner of the components of θ through the algorithm, i.e. the 3 univariate parameters of interest.

- When all components of θ are updated simultaneously, it is recommended to target an acceptance rate of around 0.20. [Roberts et al. \(1997\)](#) have shown that, under quite general conditions, the asymptotically optimal acceptance rate is 0.234. (for target density that has a symmetric product form) This quantity has been verified by [Sherlock et al. \(2009\)](#). It holds for the *Metropolis-Hastings* algorithm.
- When the components are updated one at a time, an acceptance rate of around 0.40 is recommended. It holds for the *Gibbs sampler* algorithm.

Let's (see ? for example for the first case)

Gelman-Rubin diagnostic : the \hat{R} statistic

As discussed in [item 4](#) above

Geweke diagnostic

Thinning

iteration k is stored only if $k \bmod \text{thin}$ is zero (and if k greater than or equal to the burn-in B).

This typically reduces the precision of posterior estimates, but it may represent a necessary computational saving.

5.3.2 The problem of auto and cross-correlations in the chains

There exists 2 problems of correlations in the output delivered by a MC.

- **Autocorrelation** is the
- **Cross-correlation**

5.4 Posterior Predictive

notation for posterior ? π or f

As discussed in [item 2](#) above, prediction is of important interest in EVT, and this is "facilitated" in the Bayesian paradigm. This also permits a more straightforward quantification of the inferential uncertainty associated.

Definition 5.3 (Posterior Predictive density). *Let X_{m+1} denotes a (one-step-ahead) future observation with density $f(x_{m+1}|\theta)$. Then we define the Posterior Predictive density of a future observation X_{m+1} given \mathbf{x} as*

$$\begin{aligned} f(x_{m+1}|\mathbf{x}) &= \int_{\Theta} f(x_{m+1}, \theta|\mathbf{x}) \cdot d\theta = \int_{\Theta} f(x_{m+1}|\theta) \cdot \pi(\theta|\mathbf{x}) \cdot d\theta \\ &:= \mathbb{E}_{\theta|\mathbf{x}}[f(x_{m+1}|\theta)] \end{aligned} \quad (5.5)$$

where the last line emphasizes that we can evaluate $f(x_{m+1}|\mathbf{x})$ by averaging over the different possible parameter values.

△

The uncertainty in the model is reflected here through $\pi(\theta|\mathbf{x})$ while the uncertainty due to variability in future observations is also reflected through $f(x_{m+1}|\theta)$.

Definition 5.4 (Posterior Predictive probability). *The posterior predictive probability of X_{m+1} exceeding some threshold x is accordingly given by*

$$\begin{aligned} \Pr\{X_{m+1} > x \mid \mathbf{x}\} &= \int_{\Theta} \Pr\{X_{m+1} > x \mid \theta\} \cdot \pi(\theta|\mathbf{x}) \cdot d\theta \\ &= \mathbb{E}_{\theta|\mathbf{x}}[\Pr(X_{m+1} > x \mid \theta)] \end{aligned} \quad (5.6)$$

△

This quantity is often of interest in EVT as we are rather concerned with the probability of future unknown observable exceeding some threshold.

However, this quantity is difficult to obtain analytically. Hence, we will more rely on simulated approximations. Given a sample $\theta_1, \dots, \theta_r$ from the posterior $\pi(\theta|\mathbf{x})$, we use

$$\Pr\{X_{m+1} > x \mid \mathbf{x}\} \approx r^{-1} \sum_{i=1}^r \Pr\{X_{m+1} > x \mid \theta_i\}, \quad (5.7)$$

where $\Pr\{X_{m+1} > x \mid \theta_i\}$ follows directly from $f(x|\theta)$.

We will now analyse more in-depth the numerical computations in the Bayesian paradigm or how we can get numerically a sample of the posterior distribution.

5.5 Bayesian Predictive Accuracy for Model Validation

5.5.1 Cross-validation for predictive accuracy

When having large amount of data, we can use a well-known and widely used technique coming from Machine Learning. That is, dividing the dataset between a training (typically 75% of the whole set) and a test set containing the remaining observations. For example, having N draws $\theta^{(1)}, \theta^{(2)}, \dots, \theta^{(N)}$ coming from the posterior $\pi(\theta|x_{train})$, we can score each value using (?)

$$\log \left[N^{-1} \sum_{t(i)=1}^N f(x^*|\theta^{(t)}) \right]. \quad (5.8)$$

However, we often do not have large amounts of data. Henceforth, we can use the *cross-validation* technique which is more relevant in smaller dataset, but which is computationally more demanding. There exists several variants of them.

Leave-one-out cross-validation

The *Leave-One-Out* (LOO) cross-validation is the

K -fold cross-validation

[Vehtari et al. \(2016\)](#)

Or we can use other criteria which avoid the computations. The basic approach is to use $\log f(x|\bar{\theta}) - p^*$ for N draws $\theta^{(1)}, \dots, \theta^{(N)}$ from $\pi(\theta|x)$. where p^* represents the effective number of parameters and $\bar{\theta}$ the posterior mean. Several methods exists using this idea. We will see the two most important.

Deviance Information Criterion

The *Deviance Information Criterion* (DIC) was first used by [Spiegelhalter et al. \(2002\)](#) and use the following estimate for the effective number of parameters

$$p^* = 2 \cdot \left(\log f(x|\bar{\theta}) - N^{-1} \sum_{t=1}^N \log f(x|\theta^{(t)}) \right) \quad (5.9)$$

It is defined on the deviance scale and smaller DIC values indicate better models.

$$\text{DIC} = 2 \log f(x|\bar{\theta}) - \frac{4}{N} \sum_{t=1}^N \log f(x|\theta^{(t)}) \quad (5.10)$$

Widely Applicable Information Criterion

The *Widely Applicable Information Criterion* (WAIC) is a more recent approach proposed by [Watanabe \(2010\)](#) and is given by

$$\text{WAIC} = 2 \sum_{i=1}^n \left[\log(\mathbb{E}_{\theta|x} f(x_i|\theta)) \right] - \mathbb{E}_{\theta|x} \log f(x_i|\theta) \quad (5.11)$$

or

$$\text{WAIC} = \sum_{i=1}^n \left[2 \log \left(N^{-1} \sum_{t=1}^N f(x_i|\theta^{(t)}) \right) - \frac{4}{N} \sum_{t=1}^N \log f(x_i|\theta^{(t)}) \right] \quad (5.12)$$

There exists for sure other several methods, as proposed by [Gelman et al. \(2014\)](#).

'LOO and WAIC have various advantages over simpler estimates of predictive error such as AIC and DIC but are less used in practice because they involve additional computational steps'

For each generated chains with dispersed starting values, we evaluate separately the information criteria. The discrepancies between the chains are small (?), which is a good sign.

5.6 Bayesian Inference ?

5.6.1 Bayesian Credible Intervals

The Bayesian *credible intervals* are inherently different from the frequentist's confidence intervals. In the Bayesian intervals, the bounds are treated as fixed and the estimated parameter as a random variable, while in the frequentist's setting, bound are random variables and the parameter is a fixed value.

There exist mainly two kinds of credible interval in the Bayesian sphere :

- The *Highest Probability Interval* (HPD) which is defined as the shortest interval containing $x\%$ of the posterior probability, e.g. if we want a 95% HPD interval (ξ_0, ξ_1) for ξ :

$$\int_{\xi_0}^{\xi_1} \pi(\xi|\mathbf{x})d\xi = 0.95 \quad \text{with} \quad \pi(\xi_0|\mathbf{x}) = \pi(\xi_1|\mathbf{x}). \quad (5.13)$$

It is often the preferred interval as it gives the parameter's values having the highest posterior probability.

- The Quantile-based credible intervals or *equal-tailed interval* picks an interval which ensures a probability of being below this interval as likely as of being above it. For some posterior distribution which are not symmetric, this could be misleading, thus it is not the most recommended interval. (see ..) However, these are often easily obtained when we have a random sample of the posterior...(?)

5.6.2 Distribution of Quantiles : Return Levels

The Markov chains generated can be transformed to estimate quantities of interest such as quantiles and hence return levels.

The values can be retrieved in the same manner as we have done in the GEV frequentist setting in (1.33). If the df F associated is GEV then $y_m = -\log(1 - m^{-1})$, and if F is GPD then $y_m = m^{-1}$.

r_m is the quantile corresponding to the upper tail probability $p = m^{-1}$.

We can use the values of the samples generated by the posterior to estimate features of this distribution. (... see edbayes)

5.7 Bayesian Model Averaging

Part II

Experimental Framework : Nonstationary Extreme Value Analysis of Maximum Temperatures

INTRODUCTION TO THE ANALYSIS

Contents

Repository for the code : R Package	59
Visualization Tool : Shiny Application	59
6.1 Presentation of the Analysis : Temperatures from Uccle	60
6.1.1 Open shelter vs Closed shelter	60
6.1.2 Comparisons with freely available data	60
6.2 First Analysis : Block-Maxima	61
6.2.1 Descriptive Analysis	61
6.2.2 First visualization with simple models	61
6.2.3 Deeper Trend Analysis : Splines derivatives in GAM	62
Pointwise vs Simultaneous intervals	62
Methodology	62
Final Results	63
6.3 Comments and Structure of the Analysis	64

In this part, we will focus on the application of the methods seen during the theoretical part.

Repository for the code : R Package

An **R package** has been created and can be easily downloaded from its **repository** :

<https://github.com/proto4426/PissoortThesis>

by following the instructions in the README. It follows the standard structure of an usual R package (see e.g. [Leisch \(2008\)](#)) and make use of the `roxygen2` package for the documentation.

For your (best?) convenience, an external folder has been created in this repository containing all the scripts created during this thesis. It allows to reproduce all the results, and sometimes with more details, tables or plots. It can be retrieved directly by going into the **/Scripts-R/** folder from the repository. We will notice in the beginning of each of the following chapter to which script(s) it corresponds. We also explain in more details the structure of the github repository in [appendix D](#) to give you a better idea of what each folder and each file is about.

Visualization Tool : Shiny Application

A small Shiny application has started to be developed and can be run directly from R after having the package loaded in your environment, by only typing

```
runExample() # in the R console. Then, choose the propositions displayed
```

and write it in (). So far, it contains an application based on figure 1.1 from which you can visualize the GEV distribution and the influence of the parameters ('GEV_distributions'). There are also applications dedicated to the annual maxima that can you can smoothly visualize and an application for the simulations of the GAM model with splines and the coverage visualization (section 6.2). Actually, it summarizes the contents of figure 6.1 It was difficult to do more because we could not publicly disclose the datasets provided by the IRM.

The following analysis relies on `1inter_stationary.R` and `1intro_trends(splines).R` codes from the `/Scripts-R/` folder of the github repository.

6.1 Presentation of the Analysis : Temperatures from Uccle

The data used during this thesis have been gathered directly from the "Institut Royal de Météorologie" (IRM). We were provided data...

6.1.1 Open shelter vs Closed shelter

For **meteorological considerations**, it is better to consider temperature's analysis in **closed shelters**, following for example [Lindsey and Newman \(1956\)](#). Indeed, and thanks to grateful advices from mr Tricot working for the IRM :

- It can
- It

6.1.2 Comparisons with freely available data

A similar dataset is freely available on the internet¹. It was a project initially performed by the KMNI and which was used for example in [Beirlant et al. \(2006\)](#). However, we were reticent to simply analyze these data as we know that it is hard to trust internet's data, even if they come from well-known "authorities". After having made all these comparisons analysis (see start of code...), we remark effectively that there are differences in these two datasets, and hence large errors of measures can easily occur in unofficial data. It confirms the fact that is important to get reliable data if one wants to make reliable analysis. However, these differences tend to be much smaller when considering the "open shelter" version (54% of equal measurements in closed shelter VS 14.4% in the closed case). For this reason, we have confidence that this public datasets is dealing with open shelter temperatures data.

¹<http://lstat.kuleuven.be/Wiley/Data/ecad00045TX.txt>

6.2 First Analysis : Block-Maxima

6.2.1 Descriptive Analysis

Taking yearly blocks leaves us with $n = 116$ data which seem justifiable for further GEV analysis and the conditions to hold (see)

6.2.2 First visualization with simple models

We present our series of yearly maxima in figure 6.1 where we introduce 3 **models for the trend** :

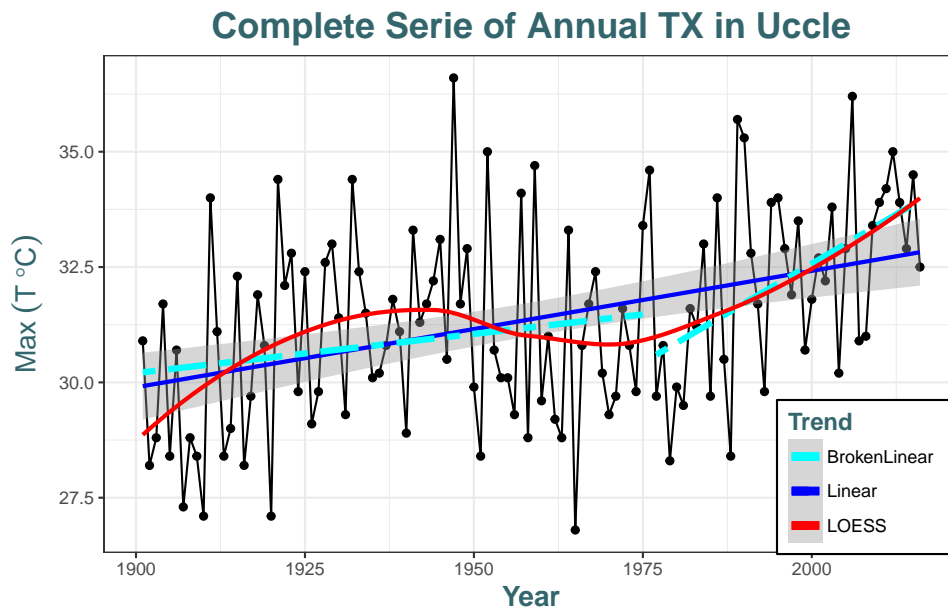


Figure 6.1: representing the **yearly maxima** together with three first models trying to represent the trend. Note that shaded grey line representing the standard errors (and not a confidence interval) of the linear trend.

- *Linear regression* (in blue) which is a **parametric** fit. We remark that it is slightly but significantly increasing over time (p-value $\approx 10^{-5}$).
- *Local Polynomial regression* or LOESS (in red) which is a **nonparametric** fit. It tells us that the yearly maxima process is rather "dispersed". For example, the drop in the series visible around years 1950 to 1975 is probably due to a 'random effect' rather than a real decrease or freezing of the maximum temperatures at this time. Moreover, it will disappear if we change the parameter controlling the degree of smoothing. We will assess that more formally in the [next section](#).
- *Broken-linear regression* (in cyan) which is a **parametric** fit. We wanted to emphasize visually the difference in trend between the period [1901-1975] and [1976-2016]. The two periods have been chosen arbitrarily and we will study that more formally in the [next section](#).

The code which provide all the tools to retrieve the presented results are left in appendix, but in the numeric version only because it is very heavy. this also enables you to get all insights (..)

6.2.3 Deeper Trend Analysis : Splines derivatives in GAM

Apart from the significant increasing linear trend, we did not find formal results. Moreover, we see from the series that a linear trend to the entire series makes little sense. Regarding the broken-linear trend for example, we would like to assess if there is indeed a difference in the slope of trend over time.

We assume the reader is at ease(?) with *Generalized Additive Models* (GAM) which have been developed by [Hastie and Tibshirani \(1986\)](#). We also assume he knows about penalized *splines* and splines smoothing, where a good introduction is found in [Ruppert et al. \(2003, section 3\)](#). This section uses the code `/Scripts-R/1intro_trends(splines).R`.

Pointwise vs Simultaneous intervals

Simultaneous confidence intervals : Following [Ruppert et al. \(2003, section 3, section 4.9, section 6.5\)](#) which uses a simulation-based approach to generate a simultaneous interval [Marra and Wood \(2012\)](#)

From the pointwise confidence intervals we can say that (example) $f(1980)$ has 95% chance to lie within $(-1,0)$ (say) and $f(2000)$ has also 95% to lie within $(0.2,1.2)$ BUT it is a fallacy to say simultaneously that both are contained in these intervals at the same time with 95% confidence. [Ruppert et al. \(2003, section 6.5\)](#)

Methodology

- We fitted a simple GAM model relying on the `mgcv` package from [Maindonald and Braun \(2006\)](#). Then, after looking at the correlation structure the normalized residuals (ACF and PACF, see figure C.2 in [appendix C](#)), we detected very slight patterns.
- As the selection from figure C.2 is not trivial, we fitted several time series models for the residuals. It is not necessary to consider too complex models so we stopped at 2 additional degrees of freedom. The results are represented in table C.1 in [appendix C](#) where we see that the BIC, which is more penalizing complex models, will prefer the independent model for the residuals but the AIC will not. Hence, we conducted likelihood ratio tests which confirmed. The diagnostics of the model in figure C.3 are good. As we took a Gaussian (identity) link, our model can be written as

$$Y_{\text{GAM}}(\text{year}) = \alpha + f_{(k)}(\text{year}) + \epsilon, \quad \epsilon \sim \text{WN}. \quad (6.1)$$

(!! WN or MA(1° for the residuals ? Depend on the analysis !!! Graph below is WN)

where f is modelled by smoothing splines. Y models the TX here. It is not important in our case to perform cross-validation to choose the dimension of the basis k and we set it to 20 to ensure a reasonable degree of smoothness.

- To account for uncertainty, we decide to simulate $M = 10^5$ draws (it is quite fast) from the posterior of the GAM model in (6.1). Hence, the confidence intervals (or bands) that we will compute will be of the form of Bayesian credible intervals that we have discussed in [section 5.6.1](#). We display 50 of the posterior draws in figure 6.2 displaying both pointwise (in yellow) and simultaneous (in red) intervals. We can already point out that pointwise intervals seem too narrow.

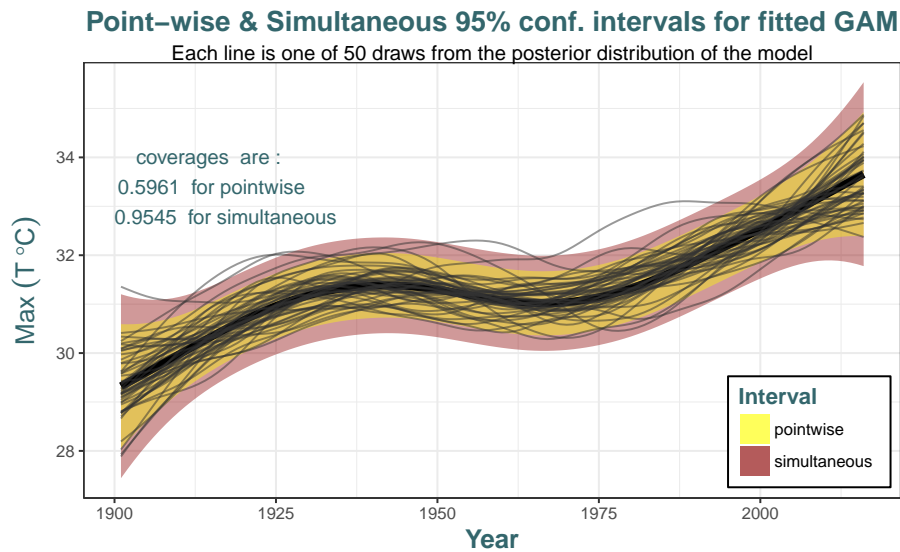


Figure 6.2: displays draws from the posterior distribution of the model. Notice the large uncertainty associated to the posterior draws (grey lines). The coverages are calculated for $M = 10^5$ simulations.

We can see the two significant increase in trend from figure C.4 let in appendix C. The decrease we have pointed out in the preceding section with LOESS is hence not significant. However, we will now see more precisely that this is not the case when doing the correction for simultaneous intervals and compare the results.

Coverage Analysis

We would like to formalize the inadequacy of the pointwise intervals. Hence, we computed posterior draws of the fitted GAM and we looked how many draws lie within each intervals.

Table 6.1: Proportion of the M posterior simulations which are covered by the confidence intervals

Coverage at 95%	$M = 20$	$M = 100$	$M = 10^3$	$M = 10^6$
<u>Pointwise</u>	40%	63%	61.1%	59.463%
<u>Simultaneous</u>	80%	91%	94.9%	95.019%

We clearly see that this indeed converges to 95% for the simultaneous interval while it converges rather to $\approx 60\%$ for the pointwise interval.

We have build a Shiny application from this graph to better visualize the impact of the number simulations on the confidence intervals and how their coverage vary, both visually and quantitatively. You can load it from the package with

```
runExample( 'splines_draws' )
```

from which you can check the results of the table (using the default seed=99)

Note that the results are similar if we do the experiment on the first derivatives $f'(\cdot)$ of the splines rather than on the splines itself

Final Results

Some remarks from the two plots in figure 6.3 :

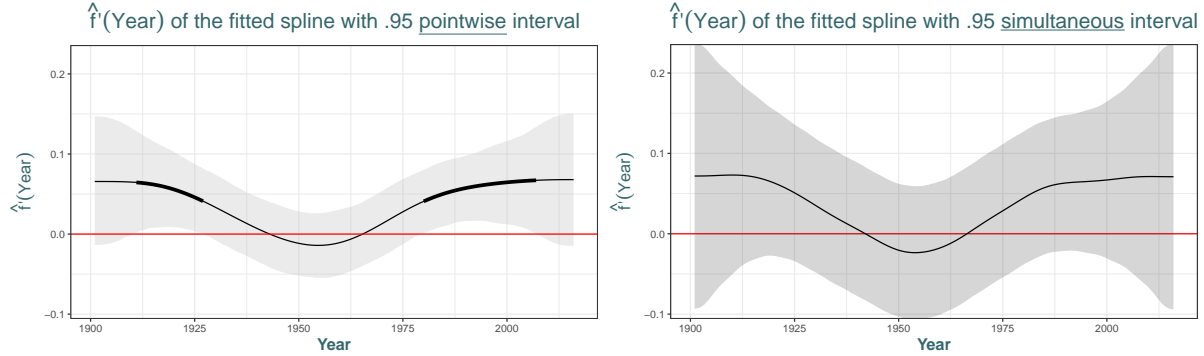


Figure 6.3: Plots of the first derivative $f'(\cdot)$ of the estimated splines on the retained GAM model. Grey area represents the 95% confidence interval. Sections of the spline where the confidence interval does not include zero are indicated by thicker sections.

- Looking together with the series in 6.1, we can make the link between the trend behaviour of the series and the splines derivative which models accurately the slope behaviour of the trend.
- Whence we can notice the increasing trend but with decreasing slope until ≈ 1945 where it becomes is decreasing. Then the slope increase and the trend starts to re-increase around 1962. This upward trend that the series of annual maxima is facing for this last period brings light to the climate warming we are all talking about. However,
- Whereas the pointwise confidence interval include significant regions (i.e. when 0 is not included in the interval), the simultaneous interval which accounting for the increase in uncertainty, have no more significant regions. We can conclude that the

We can see that only the 2 increasing periods from the start and at the end of the series are significant, while the decreasing period from... (see the red line) is more likely to be the subject of randomness. Moreover, we remark that the

6.3 Comments and Structure of the Analysis

We have remarked that there is indeed an upward (linear) trend for the series of yearly maxima but which is not significant when doing the correction for simultaneous intervals. We have then remarked that the **nonstationary** component is present. After having... we will now go through with the more specific subjec of this thesis, that is the extreme value analysis. Hence, after having presented a stationary analysis in section 7, we will make an in-depth nonstationary analysis in section 8.

POT and GEV : As you have seen, the analysis in POT or in GEV involves different methods and different data. For the kind of this text, we decide not to display the results of the POT analysis to keep the text not too enormous. Henceforth, we will be able to focus on the GEV analysis. But **note that the**

analysis in POT have been done and are available on the same [github repository](#) presented above. [Appendix D](#) summarize its contents.

Results are not shown here but this is (visually) less pronounced for minima. We have also conducted some analysis by dividing the dataset by seasons, by hot and cold months (July-August, January-February), ... Comparisons are interesting (and are also available on the repository) but this is not at the core of this text. We will then only **focus on a GEV analysis on yearly maxima**.

FIRST ANALYSIS BY GEV

Contents

R packages for EVT	66
7.1 Estimation of the Model	67
7.1.1 Maximum Likelihood	67
7.1.2 Other Methods	68
7.2 Diagnostics	68
7.3 Return Levels	68
7.3.1 Profile Likelihood	68
7.4 Comments and Comparisons with POT	68

PUT the examples right in the place where it is mentioned in the theory! "As we have seen in section 2.1.1.... and in section 2.2.2....."

In this first analysis, we rely on

This analysis relies on `1inter_stationary.R` code from the **/Scripts-R/** folder of the github repository.

The block-length selection an important issue of the analysis. It is important to choose a block-length which is large enough for the limiting arguments supporting the GEV approximation (see (1.8)) to be valid, either a large bias in the estimates could occur. For example, if this is too short, the maxima may be too close of each other to assume independence. But a large block-length implies less data to work with, and thus a large variance of the estimates. A compromise must be found between bias and variance and as pointed out in the previous section, yearly blocks seem justifiable for both this reason but also for their interpretability.

R packages in EVT

A plenty of packages exist for modelling extreme values in R. We have explored some of them and we used some of them. For "basic" EVT analysis, we must name the following :

- **ismev**, **evd**, **extRemes** (good for a wide nonstationary analysis with POT and nice tutorials, see e.g. [Gilleland and Katz \(2016\)](#)) , **POT** (see [Ribatet \(2006\)](#)), **evir**, **fEx-tremes**, ...

Whereas lots of the package are doing the same analysis but with different tools, we relied mostly on `ismev` as it is the package used in the book of [Coles \(2001\)](#).

7.1 Estimation of the Model

Whereas the whole content of chapter 1 is important to understand the concepts used in this section, we will now be based on inferential methods discussed in section 1.7

7.1.1 Maximum Likelihood

Relying on the packages named above but also by checking it manually, i.e. by numerically solving the optimization problem, that is minimizing the negative log-likelihood, with the `nlm` routine using a Newton-Raphson algorithm (originated from Dennis and Schnabel (1987)). This is based on approximating the log-likelihood by a quadratic function, the second order Taylor series approximation of the log-likelihood for a given point. The results are shown in table 7.1 :

Table 7.1: shows the maximum likelihood estimation of the three GEV parameters

	Location μ	Scale σ	Shape ξ
Estimates (s.e.)	30.587 (0.216)	2.081 (0.155)	−0.254 (0.067)

From this table 7.1, an important thing to note is the value of the **shape** parameter which is negative which means that we are under a Weibull-type distribution. From figure ?? this means that the density distribution has the form of the red line, i.e. having a . Moreover, we confirmed that by doing a likelihood ratio test comparing this distribution with a Gumbel distribution. We obtained a p-value of 10^{-3} leading to a rejection of the Gumbel hypothesis. This obviously implies rejection of the Fréchet distribution.

The Weibull-type implies that the distribution have an estimated right endpoint given by $\hat{x}_* = \hat{\mu} - \hat{\sigma} \cdot \hat{\xi}^{-1} = 38.77$. Comparing this value with the maximum value of the series (=36.6) tells us that the sample properties of this model take into account the uncertainty, from the fact that there are only 116 years of data. Hence, it allows to go beyond this maximum value, with very small probability. This is also highlighted in figure 7.2 (right plot) where we remark that there are still probability mass beyond the minimum and the maximum values of the series.

The **profile** log-likelihood intervals for the parameters are shown in figure C.5 in appendix C, provided by the `ismev` package. These intervals are constructed in the following way : Search for the horizontal line Subtracting the maximum log-likelihood ℓ half the corresponding upper quantile of the χ^2_{df} for the $df = 1$ parameter of interest. We notice that :

- Even at 99%, the interval for $\hat{\xi}$ does not contain 0 supporting our statement that the distribution is left heavy-tailed and right bounded.
- The intervals do not present much asymetries. In fact, this will be more relevant for return levels as we will see in section

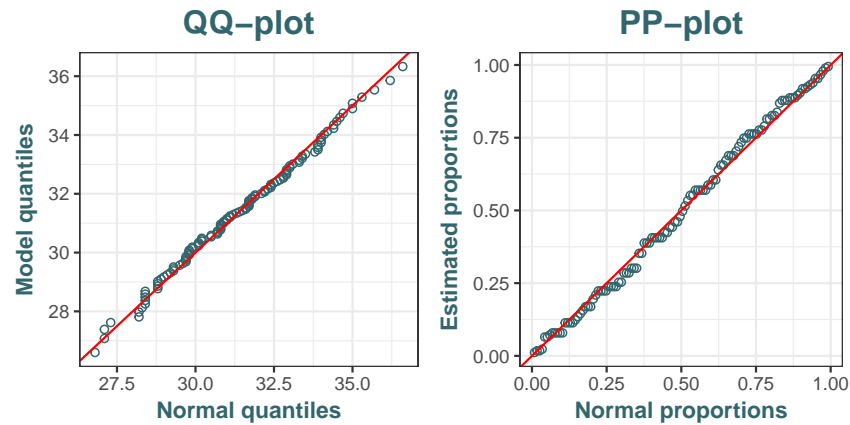


Figure 7.1

7.1.2 Other Methods

7.2 Diagnostics

7.3 Return Levels

Explore why the return levels go beyond the right endpoint of the distribution (when $\xi < 0$ as here), for which return period, etC...

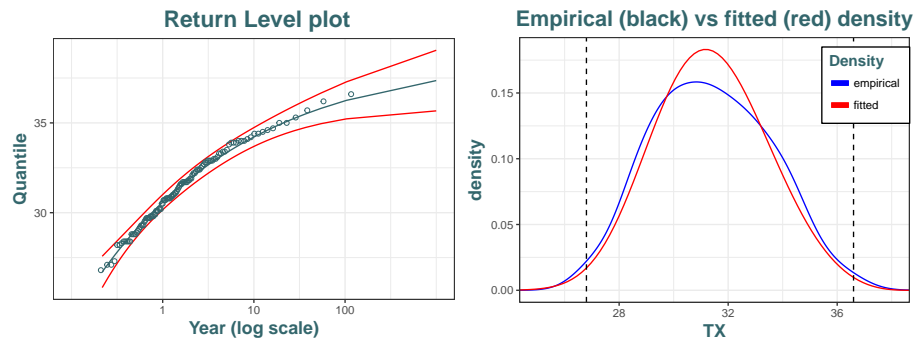


Figure 7.2: The left and right vertical dotted lines represent respectively the minimum and the maximum value of the yearly maxima series.

7.3.1 Profile Likelihood

7.4 Comments and Comparisons with POT

STATIONARY AND NONSTATIONARY GEV ANALYSIS

Contents

8.1 Stationary Analysis	69
8.2 Nonstationary Analysis	69
8.3 Improvements with Neural Networks	69

This analysis relies on `2Nonstationary.R` and `2NeuralSNets.R` code from the `/Scripts-R/` folder of the github repository.

TABLE with nested models (gumbel, GEV, + linear trend, etc etc)

8.1 Stationary Analysis

In [section 3.1](#) we have proven that

8.2 Nonstationary Analysis

8.3 Improvements with Neural Networks

BAYESIAN ANALYSIS

Contents

9.1 From <code>evdbayes</code> R package : MH algorithm	70
9.2 From Our Functions (R package)	70
9.3 From HMC algorithm using STAN language	70
9.4 Comparisons	71
9.4.1 STAN	71
9.5 Comparison with frequentists results	71

This analysis relies on all the codes which filenames start by "Bayes" from the `/Scripts-R/` folder of the github repository. All the functions created that are used and are also made available through the package are in `/R/BayesFunc.R`.

"It is often the case that more than one model provides an adequate fit to the data. Sensitivity analysis determines by what extent posterior inferences change when alternative models are used" book risk analysis other section pp.2.

"The basic method of sensitivity analysis is to fit several models to the same problem. Posterior inferences from each model can then be compared."

9.1 From `evdbayes` R package : MH algorithm

The `evdbayes` is the only package available on CRAN for (see [Ribatet \(2006\)](#)) is a very old package which has not been updated since a while. We had problems to understand both its structure and the famous "black-box"

9.2 From Our Functions (R package)

9.3 From HMC algorithm using STAN language

The problem is maybe from [1.7.1](#). The parameter ξ is relatively near the region that could be problematic, causing convergence issues.

9.4 Comparisons

Hartmann and Ehlers (2016) We can calculate the effective sample size (ESS) using the posterior samples for each parameter :

$$\text{ESS} = N \cdot \left(1 + 2 \sum_k \gamma(k)\right)^{-1} \quad (9.1)$$

where N is still the number of posterior samples and $\gamma(k)$ are the monotone lag k sample autocorrelations. We can thus interpret this as the number of effectively independent samples.

9.4.1 STAN

Benefits :

- Allows more flexibility (?) through the mathematical formulation of the formula
- It is really smoother and clearer (straightforward) for this kind of problems

Drawbacks :

- New language with all the problems/errors arising when learning it.

9.5 Comparison with frequentists results

CONCLUSION

During this thesis, we have statistically assessed the presence of a trend in the extreme temperatures in Uccle. We first detected that the trend is significative by the method of linear regression. We also discovered that the best fitted GEV model is the one with a linear trend in the location parameter.

"A key issue in applications is that inferences may be required well beyond the observed tail of the data, and so an assumption of stability is required:" [Davison et al. \(2012\)](#)

"Another approach would be to use something other than time as the covariate in the model. For instance, one could imagine linking temperature data directly to CO2 level rather than time. However, linking to a climatological covariate makes extrapolation into the future more difficult, as one would need to extrapolate the covariate as well. No obvious climatological covariate comes to mind for the Red River application. "

Timescale-uncertainty effects on extreme value analyses seem not to have been studied yet. For stationary models (Sect. 6.2), we anticipate sizable effects on block extremes–GEV estimates only when the uncertainties distort strongly the blocking procedure. For nonstationary models (Sect. 6.3), one may augment confidence band construction by inserting a timescale simulation step (after Step 4 in Algorithm 6.1) [Mudelsee \(2014, pp.262\)](#)

!!!! not put too much references in the text !!!!!

Managing references :

(!!!! delete unuseful equation numbering ?!!!!)

finalize hyperref in the text. BE CAREFULL of the numbering written and the real ("hyperrefed" numbering section) +finalize also for def, thm, ... (?)

Replace "distribution functions" by "cdf",.... + all other abbreviations (et, gev,...) -> put it then in the list of abbr.

Careful with notation with X or Z in the (c)pdf, likelihood, RV,...

ecrire correctement les referencesz a la fin (unsrt85.bst, ...) enlever trop de noms, champs inutiles,...

voir notation vectors (en gras ou avec bar en bas?)

attention aux notation homogenes (ex : partie stationnary,...) -> pour denoter les maximums,....

notations for sequences !! attention aux "iid", "n random variables",...

Delete page counter for part pages,...

be prudent that all "methods" are listed (numbered) in (each) ToC

add square brackets [] for cite ?

Careful for bold symbols, especially in bayesian

PissoortThesis:: Behind our functions

[MOTS DEXPLICATIONS SUR TTES LES FORMULES (domain attraction condition, etc...)]

Finish (Bayesian) documentation in the package.

Appendix

APPENDIX A

STATISTICAL TOOLS FOR EXTREME VALUE THEORY

A.1 Tails of the distributions

Heavy-tailed

The distribution of a random variable X with distribution function F is said to have a *heavy right tail* if

$$\lim_{n \rightarrow \infty} e^{\lambda x} \Pr\{X > x\} = \lim_{n \rightarrow \infty} e^{\lambda x} \bar{F}(x) = \infty, \quad \forall \lambda > 0. \quad (\text{A.1})$$

More generally, we can say that a random variable X has heavy tails if $\Pr\{|X| > x\} \rightarrow 0$ at a polynomial rate. In this case, note that some of the moments will be undefined.

Fat-tailed The distribution of a random variable X is said to have a *fat tail* if

$$\lim_{x \rightarrow \infty} \Pr\{X > x\} = x^{-\alpha}. \quad (\text{A.2})$$

Long-tailed The distribution of a random variable X with distribution function F is said to have a *long right tail* if $\forall t > 0$,

$$\lim_{x \rightarrow \infty} \Pr\{X > x + t | X > x\} = 1 \Leftrightarrow \bar{F}(x + t) \sim \bar{F}(x) \quad \text{as } x \rightarrow \infty. \quad (\text{A.3})$$

Light-tailed

Conversely, we say that X has *light tails* or *exponential tails* if its tails decay at an exponential rate, i.e.

$$\lim_{x \rightarrow \infty} \Pr\{|X| > x\} = e^{-x} \quad (\text{A.4})$$

An intuitive example of a distribution with exponential tails

A.2 Convergence concepts

Convergence in distribution

We say that a sequence X_n with df F_n converges in distribution to X with df F , if

$$F_n(x) := \Pr\{X_n \leq x\} \longrightarrow \Pr\{X \leq x\} := F(x), \quad (\text{A.5})$$

at all continuity points of F . It means that, for large n , $\Pr\{X_n \leq t\} \approx \Pr\{X \leq t\}$. We denote this by $X_n \xrightarrow{d} X$.

Convergence in probability

We say that a sequence X_n converges to X in probability if, $\forall \epsilon > 0$,

$$\Pr\{|X_n - X| > \epsilon\} \rightarrow 0, \quad n \rightarrow \infty. \quad (\text{A.6})$$

Hence, it means that the probability of the difference between X_n and X goes to 0 as n is large. We denote this by $X_n \xrightarrow{P} X$.

An example of application of this convergence is the *Weak Law of Large Numbers*.

Definition A.1 (Weak Law of Large Numbers). *Let a sequence of R.V. $\{X_i\}_{iid}$ be defined of the same probability space with mean μ and variance $\sigma^2 < \infty$. Then, we know that the difference between \bar{X}_n and μ will go to 0 in probability, i.e. $\bar{X}_n \xrightarrow{P} \mu$. \triangle*

But this law actually makes a stronger convergence, following [Kolmogorov et al. \(1956\)](#)

Almost Sure Convergence

This is the type of stochastic convergence that is most similar to pointwise convergence known from elementary real analysis.

We say that a sequence of random variables X_n converges almost surely (or with probability one) to X if

$$\Pr\{X_n = X\} = 1, \quad n \rightarrow \infty. \quad (\text{A.7})$$

We can denote this by $X_n \xrightarrow{\text{a.s.}} X$. This means that the values of X_n approach the value of X , in the sense that events for which X_n does not converge to X have probability 0.

Well other forms of convergence do exist, but these ones are the most important in regard to EVT. However, the reader may refer e.g. to ? for more in-depth results.

A.3 Varying functions

Definition A.2 (Regularly varying function). *Let's consider the survival \bar{F} . We say that this survival function \bar{F} is **regularly varying** with index $-\alpha$ if*

$$\lim_{x \rightarrow \infty} \frac{\bar{F}(xt)}{\bar{F}(x)} = t^{-\alpha}, \quad t > 0. \quad (\text{A.8})$$

We write it $\bar{F} \in R_{-\alpha}$. △

Definition A.3 (Slowly varying function). *We say that a function f is **slowly varying** with index $-\alpha$ if*

$$\lim_{x \rightarrow \infty} \frac{f(tx)}{f(x)} = 1, \quad t > 0. \quad (\text{A.9})$$

△

We remark that a slowly varying function is a regularly varying function with index 0.

A.4 Diagnostic Plots : Quantile and Probability Plots

From [Beirlant et al. \(1996, pp.18-36\)](#), together with the nice view of [Coles \(2001, pp.36-37\)](#), we present two major diagnostic tools which aim at assessing the fit of a particular model (or distribution) against the real distribution coming from the data used to construct the model. These are called the *quantile-quantile plot* (or *qq-plot*) and the *probability plot* (or *pp-plot*).

These diagnostics are popular by their easy interpretation and by the fact that they can both have graphical (i.e. subjective, qualitative, quick) view but also a more precise (i.e. objective, quantitative, rigorous) analysis can be derived, for example from the theory of linear regression.

For these two diagnostic tools, we use the order statistics as seen (1.1) but now we rather consider an **ordered sample** of independent **observations** :

$$x_{(1)} \leq x_{(2)} \leq \dots \leq x_{(n)} \quad (\text{A.10})$$

coming from a population from which we fit the estimated model (distribution) \hat{F} and where $x_{(1)}$ (resp. $x_{(n)}$) is thus the minimum (resp. maximum) observation in the sample. We also define the **empirical distribution function**

$$\tilde{F}(x) = \frac{i}{n+1}, \quad x_{(i)} \leq x \leq x_{(i+1)}. \quad (\text{A.11})$$

\tilde{F} is an estimate of the true distribution F and hence, by comparing \hat{F} and \tilde{F} , it will help us to know if the fitted model \hat{F} is reasonable for the data.

Quantile plot

Given a ordered sample as in (A.10), a *qq-plot* consists of the locus of points

$$\left\{ \left(\hat{F}^{\leftarrow} \left(\frac{i}{n+1} \right), x_{(i)} \right) : i = 1, \dots, n \right\}. \quad (\text{A.12})$$

This graphic compares the ordered quantiles $\hat{F}^{\leftarrow} \left(\frac{i}{n+1} \right)$ of the fitted model \hat{F} against the ordered observed quantiles, i.e. the ordered sample from (A.10). We used the continuity correction $\frac{i}{n+1}$ to prevent problems at the borders. Note that a disadvantage of Q-Q plots is that the shape of the selected parametric distribution is no longer visible [Beirlant et al. \(2006\)](#)[pp.63]

Probability plot

Given the same sample in (A.10), a *probability plot* consists of the locus of points

$$\left\{ \left(\hat{F}(x_{(i)}), \frac{i}{n+1} \right) : i = 1, \dots, n \right\}. \quad (\text{A.13})$$

This graph compares the estimated probability of the ordered values $x_{(i)}$, thus from the fitted model \hat{F} , against the probability coming from the empirical distribution as in (A.11).

From these two graphical diagnostic tools, the interpretation is the same and we will consider that \hat{F} fits well the data if the plot looks linear, i.e. the points of the plots lie close to the unit diagonal.

Besides the fact that the probability and the quantile plots contain the same information, they are expressed in a different scale. That is, after changing the scale to probabilities or quantiles (with probability or quantile transforms), one can gain a better perception and both visualizations can sometimes lead contradictory conclusions, especially in the graphical inspection. Using both is thus preferable to make our model's diagnostic more robust.

APPENDIX B

BAYESIAN METHODS

B.1 Algorithms

B.1.1 Metropolis–Hastings Algorithm

The *Metropolis–Hastings* algorithm is one of the first and of the pioneering algorithm discovered by [Hastings \(1970\)](#) to compute MCMC for Bayesian analysis.

Algorithm 1: The Metropolis–Hastings Algorithm

1. Pick a starting point θ_0 and fix some number N of simulations.

2. **For** $t = 1, \dots, N$ **do**

(a) Sample proposal θ_* from a proposal density $p_t(\theta_*|\theta_{t-1})$,

(b) Compute the ratio

$$r = \frac{\pi(\theta_*|\mathbf{x}) \cdot p_t(\theta_{t-1}|\theta_*)}{\pi(\theta_{t-1}|\mathbf{x}) \cdot p_t(\theta_*|\theta_{t-1})} = \frac{\pi(\theta_*) \cdot \pi(\mathbf{x}|\theta_*) \cdot p_t(\theta_{t-1}|\theta_*)}{\pi(\theta_{t-1}) \cdot \pi(\mathbf{x}|\theta_{t-1}) \cdot p_t(\theta_*|\theta_{t-1})}.$$

(c) Set

$$\theta_t = \begin{cases} \theta_* & \text{with probability } \alpha = \min(r, 1); \\ \theta_{t-1} & \text{otherwise.} \end{cases}$$

This algorithm remains valid when π is only proportional to a target density function and thus it can be used to approximate [5.1](#).

Note that the proposal density is often chosen to be symmetric so that we will just sample under a "simple Metropolis" algorithm where r is thus simplified to be only the ratio of the posterior densities, $r = \frac{\pi(\theta_*|\mathbf{x})}{\pi(\theta_{t-1}|\mathbf{x})}$.

We can shortly summarize the *pros* and *cons* this algorithm :

- *PROS* : Very easy to program and works even for relatively complex densities.
- *CONS* : Can be very inefficient, in the sense that it will require lots of iterations before the stationary target distribution will be reached. This requires some tuning to the algorithm through

B.1.2 Gibbs Sampler

The *Gibbs Sampler* can be seen as a special case of the Metropolis-Hastings algorithm. Suppose our parameter vector θ can be divided into d subvectors $(\theta_1, \dots, \theta_d)$, and let's say in our case that each of these "subvectors" represent a single parameter, thus typically one of the three (μ, σ, ξ) , for the simplest case so that $d = 3$ in this model. At each $t = 1, \dots, N$, the Gibbs sampler samples the subvectors $\theta_t^{(j)}$ conditional on both the data \mathbf{x} and the remaining subvectors $\theta_{t-1}^{(-j)}$ at their current values. Therefore, we have $\theta_{t-1}^{(-j)} = (\theta_t^{(1)}, \dots, \theta_t^{(j-1)}, \theta_{t-1}^{(j+1)}, \dots, \theta_{t-1}^{(d)})$ and each $\theta_t^{(j)}$ is sampled from $\pi(\theta^{(j)} | \theta_{t-1}^{(-j)}, \mathbf{x})$.

Algorithm 2: PSEUDOCODE of the Gibbs Sampler

1. Pick a starting point θ_0 and fix some number N of simulations.

2. **For** $t = 1, \dots, N$ **do**
 For $j = 1, \dots, d$ **do**

 (a) Sample proposal θ_* from a proposal density $p_{t,j}(\theta_*^{(j)} | \theta_{t-1}^{(j)})$,

 (b) Compute the ratio

$$\begin{aligned} r &= \frac{\pi(\theta_*^{(j)} | \theta_{t-1}^{(-j)}, \mathbf{x}) \cdot p_{t,j}(\theta_{t-1}^{(j)} | \theta_*^{(j)})}{\pi(\theta_{t-1}^{(j)} | \theta_{t-1}^{(-j)}, \mathbf{x}) \cdot p_{t,j}(\theta_*^{(j)} | \theta_{t-1}^{(j)})} \\ &= \frac{\pi(\theta_t^{(1)}, \dots, \theta_t^{(j-1)}, \theta_*^{(j)}, \theta_{t-1}^{(j+1)}, \dots, \theta_{t-1}^{(d)} | \mathbf{x}) \cdot p_{t,j}(\theta_{t-1}^{(j)} | \theta_*^{(j)})}{\pi(\theta_t^{(1)}, \dots, \theta_t^{(j-1)}, \theta_{t-1}^{(j)}, \theta_{t-1}^{(j+1)}, \dots, \theta_{t-1}^{(d)} | \mathbf{x}) \cdot p_{t,j}(\theta_*^{(j)} | \theta_{t-1}^{(j)})}, \end{aligned}$$

 (c) Set

$$\theta_t^{(j)} = \begin{cases} \theta_*^{(j)} & \text{with probability } \alpha = \min(r, 1); \\ \theta_{t-1}^{(j)} & \text{otherwise.} \end{cases}$$

(better signs for conditional bar "|" !!!!!!!)

This algorithm depends on being able to simulate from $\pi(\theta^{(j)} | \theta_{t-1}^{(-j)}, \mathbf{x})$ which is often impossible. However, one can use Metropolis-Hastings to $\pi(\theta^{(j)} | \theta_{t-1}^{(-j)}, \mathbf{x})$, giving the above.

A special case arise if one can simulate directly so that $r = 1$, i.e. we take $p_{t,j}(\theta_*^{(j)} | \theta_{t-1}^{(j)}) = \pi(\theta_*^{(j)} | \theta_{t-1}^{(-j)}, \mathbf{x})$. (The proposal $p_{t,j}(\cdot)$ is also often symmetric, i.e. $p_{t,j}(\theta_*^{(j)} | \theta_{t-1}^{(j)}) = p_{t,j}(\theta_{t-1}^{(j)} | \theta_*^{(j)})$. But it cannot be simplified in the equation.)

It is important for our tasks to tune the average probability of acceptance to be roughly between 0.4 and 0.5 (see e.g. [Gelman et al. \(2013, chapter 11\)](#)) so that the so-generated markov-chain has desirable properties. This is done by setting the standard deviation $\sigma(j)$ of the univariate normal distribution taken for $p_{t,j}(\theta_*^{(j)} | \theta_{t-1}^{(j)})$. Whereas our $\theta^{(j)}$ are (often?) univariate, it is difficult to set each $\sigma^{(j)}$ to achieve average acceptance probabilities for all parameters. We will then use a trial-and-error approach.

Note also the increase of complexity with this sampler compared to the Metropolis-Hastings, where the nested loop implies that there are d iterations with each simulation.

pros and cons :

- *PROS* : Easy to program and, for some problems, it can also be very efficient. It is a pleasant way to split multidimensional problems into simpler (typically univariate) densities.
- *CONS* : Sometimes hard to compute analytically the conditional distributions. Not all densities can be split into pleasant conditionals equations.

Metropolis-within-Gibbs?

B.1.3 Hamiltonian Monte Carlo

<http://deeplearning.net/tutorial/hmc.html#hmc>

A difficulty we have faced and that we would like to point out for GEV models is (see p.316-317 STAN manual) t

"Most of the computation [in Stan] is done using Hamiltonian Monte Carlo. HMC requires some tuning, so Matt Hoffman up and wrote a new algorithm, Nuts (the "No-U-Turn Sampler") which optimizes HMC adaptively. In many settings, Nuts is actually more computationally efficient than the optimal static HMC! "

The *Hamiltonian Monte Carlo* (HMC)

"The Hamiltonian Monte Carlo algorithm starts at a specified initial set of parameters θ ; in Stan, this value is either user-specified or generated randomly. Then, for a given number of iterations, a new momentum vector is sampled and the current value of the parameter θ is updated using the leapfrog integrator with discretization time and number of steps L according to the Hamiltonian dynamics. Then a Metropolis acceptance step is applied, and a decision is made whether to update to the new state (θ^* ; ρ^*) or keep the existing state" ?

? have demonstrated in similar application that HMC (and Riemann manifold HMC) are much more computationally efficient than traditional MCMC algorithms such as MH.

Definition B.1 (Total energy of a closed system : Hamiltonian function). *For a certain particle; Let $\pi(\theta)$ be the posterior distribution and let $\mathbf{p} \in \mathbb{R}^d$ denote a vector of auxiliary parameters independent of θ and distributed as $\mathbf{p} \sim N(\mathbf{0}, \mathbf{M})$. We can interpret θ as the position of the particle and $-\log \pi(\theta|x)$ describes its potential energy while \mathbf{p} is the momentum with kinetic energy $\mathbf{p}'\mathbf{M}^{-1}\mathbf{p} \cdot 2^{-1}$. Then the total energy of a closed system is the Hamiltonian function*

$$\mathcal{H}(\theta, \mathbf{p}) = -\mathcal{L}(\theta) + \mathbf{p}'\mathbf{M}^{-1}\mathbf{p} \cdot 2^{-1}, \quad \text{where} \quad \mathcal{L}(\theta) = \log \pi(\theta). \quad (\text{B.1})$$

△

We define $\mathcal{X} = (\theta, \mathbf{p})$ as the combined state of the particle.

The unnormalized joint density of (θ, \mathbf{p}) is

$$f(\theta, \mathbf{p}) \propto \pi(\theta) \cdot \exp\{-\mathbf{p}'\mathbf{M}^{-1}\mathbf{p} \cdot 2^{-1}\} \propto \exp\{-\mathcal{H}(\theta, \mathbf{p})\}. \quad (\text{B.2})$$

Following [Hartmann and Ehlers \(2016\)](#), the idea is to use the Hamiltonian dynamics equations (not shown here for..) to model the evolution of a particle that keep the total energy constant. Introducing the auxiliary variables \mathbf{p} and using the gradients (..) will lead to a more efficient exploration of the parameter space

These differential equations cannot be solved so numerical integrators are required, for instance the "Störmer-Verlet" from ? which will introduce discretization. A MH step is then required to correct the error and ensure convergence. The new proposal $\mathcal{X}_* = (\theta_*, \mathbf{p}_*)$ will be accepted with probability

$$\alpha(\mathcal{X}, \mathcal{X}_*) = \min \left[\frac{f(\theta_*, \mathbf{p}_*)}{f(\theta, \mathbf{p})}, 1 \right] = \min \left[\exp \{ \mathcal{H}(\theta, \mathbf{p}) - \mathcal{H}(\theta_*, \mathbf{p}_*) \}, 1 \right]. \quad (\text{B.3})$$

As \mathbf{M} is symmetric positive definite, $\mathbf{M} = m\mathbf{I}_d$. Then we can summarize the [HMC algorithm](#) in the following, in its 'simplest' form :

$$\text{marie est la best : } \text{Moyenne} = \frac{1}{n} \sum_{i=1}^n X_i$$

Algorithm 3: The Hamiltonian Monte Carlo algorithm

1. Pick a starting point θ_0 and set $i = 1$.
 2. **Until** convergence has been reached **do**
 - (a) Sample $\mathbf{p}_* \sim N_d(\mathbf{0}, \mathbf{I}_d)$ and $u \sim U(0, 1)$,
 - (b) Set $(\theta_I, \mathbf{p}_I) = (\theta_{i-1}, \mathbf{p}_*)$ and $\mathcal{H}_0 = \mathcal{H}(\theta_I, \mathbf{p}_I)$,
 - (c) **repeat** L times
 - $\triangleright \mathbf{p}_* = \mathbf{p}_* + \frac{\epsilon}{2} \nabla_{\theta} \mathcal{L}(\theta_{i-1})$
 - $\triangleright \theta_{i-1} = \theta_{i-1} + \epsilon \cdot \mathbf{p}_*$
 - $\triangleright \mathbf{p}_* = \mathbf{p}_* + \frac{\epsilon}{2} \nabla_{\theta} \mathcal{L}(\theta_{i-1})$,
 - (d) Set $(\theta_L, \mathbf{p}_L) = (\theta_{i-1}, \mathbf{p}_*)$ and $\mathcal{H}^{(1)} = \mathcal{H}(\theta_L, \mathbf{p}_L)$,
 - (e) Compute $\alpha \left[(\theta_I, \mathbf{p}_I), (\theta_L, \mathbf{p}_L) \right] = \min \left[\exp \{ \mathcal{H}^{(0)} - \mathcal{H}^{(1)} \}, 1 \right]$,
 - (f) **If** $\alpha \left[(\theta_I, \mathbf{p}_I), (\theta_L, \mathbf{p}_L) \right] > u$ **then** set $\theta_i = \theta_L$
else set $\theta_i = \theta_I$,
 - (g) Increment $i = i + 1$ and return to [step \(a\)](#).
-

As you can see, it is not trivial. The basic idea to keep in mind is that jumping rules are much more efficient than for traditional algorithms because they learn from the gradient of the log posterior density, so they know better where to jump to. As a result, it can be MUCH more efficient.

Chains are expected to reach stationarity faster as it proposes moves to regions of higher probabilities.

pros and cons :

- *PROS* : Easy to program as we just have to write down the model. Very efficient in general, and works for all types of problems.

-
- *CONS* : Need to learn how to use STAN, less control over the sampler bu maybe it is for the best?

OTHER FIGURES AND TABLES

C.1 GEV : Influence of the Parameters on the shape of the distribution

Regarding our future application, that is maximum temperatures, it is relevant to consider values of the location parameter μ around 30 degrees.

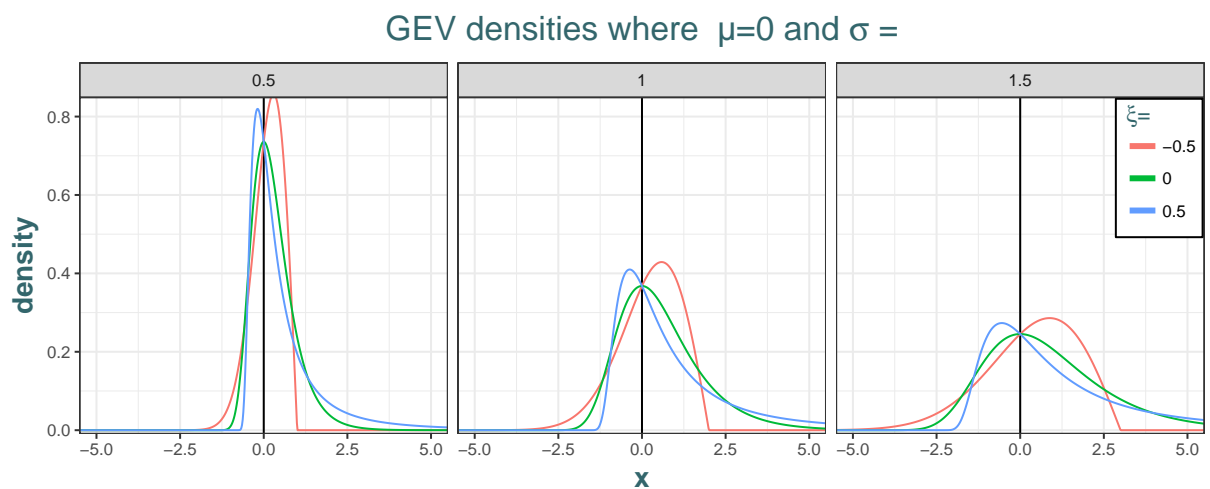
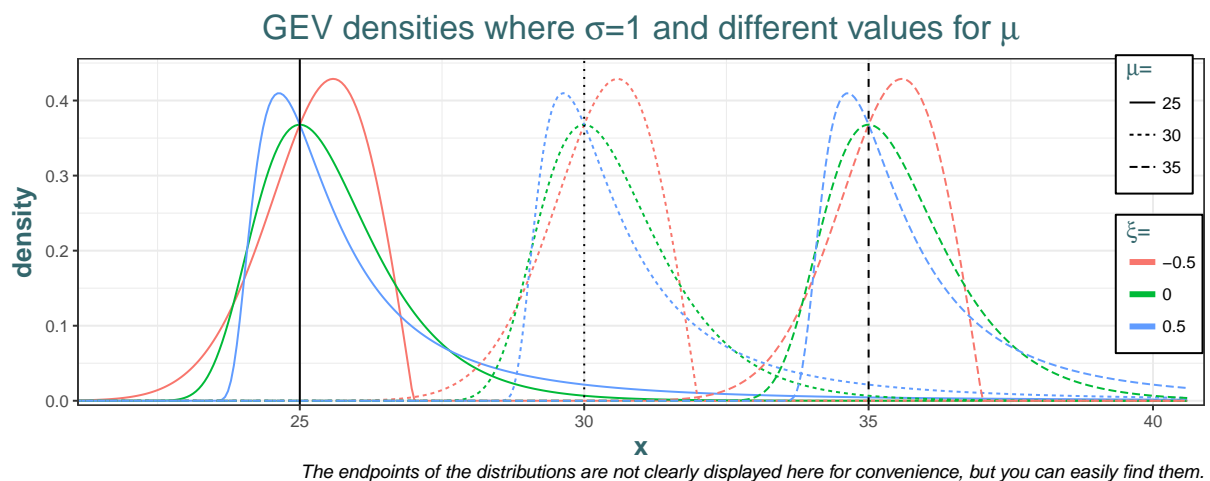


Figure C.1: GEV distribution for different values of the three parameters

C.2 Introduction of the Practical Analysis (section 6)

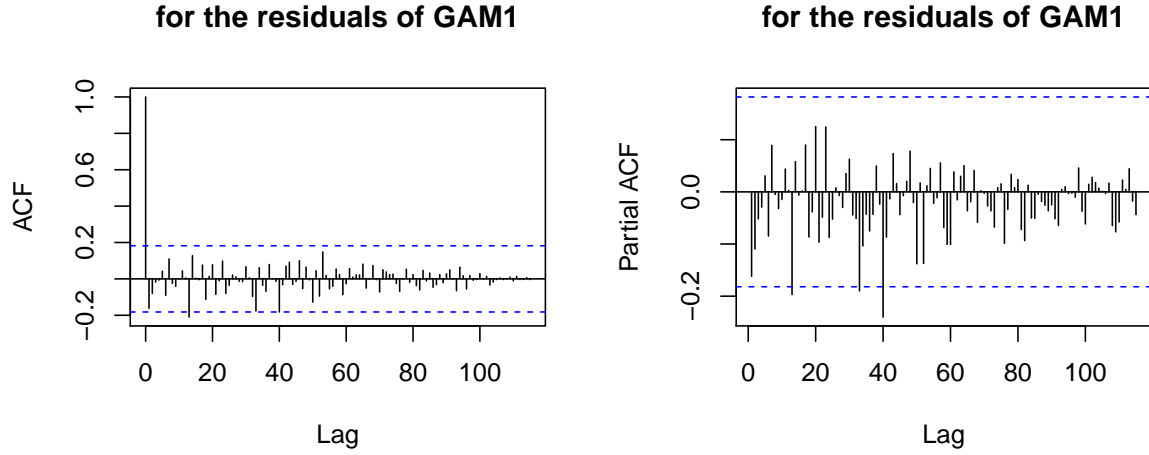


Figure C.2: ACF and PACF for the residuals of the fitted GAM model with assumed independent errors

Table C.1: compares models for the residuals of the GAM model based on AIC and BIC criterion. These criterion take into account the quality of fit (based on likelihood) but also a penalty term to penalize more complex models.

	df	AIC	BIC
Uncorrelated	4	494.635	505.650
AR(1)	5	494.356	508.124
MA(1)	5	493.706	507.474
ARMA(1,1)	6	492.511	509.033
AR(2)	6	495.133	511.654
MA(2)	6	494.698	511.219

C.3 Analysis by GEV

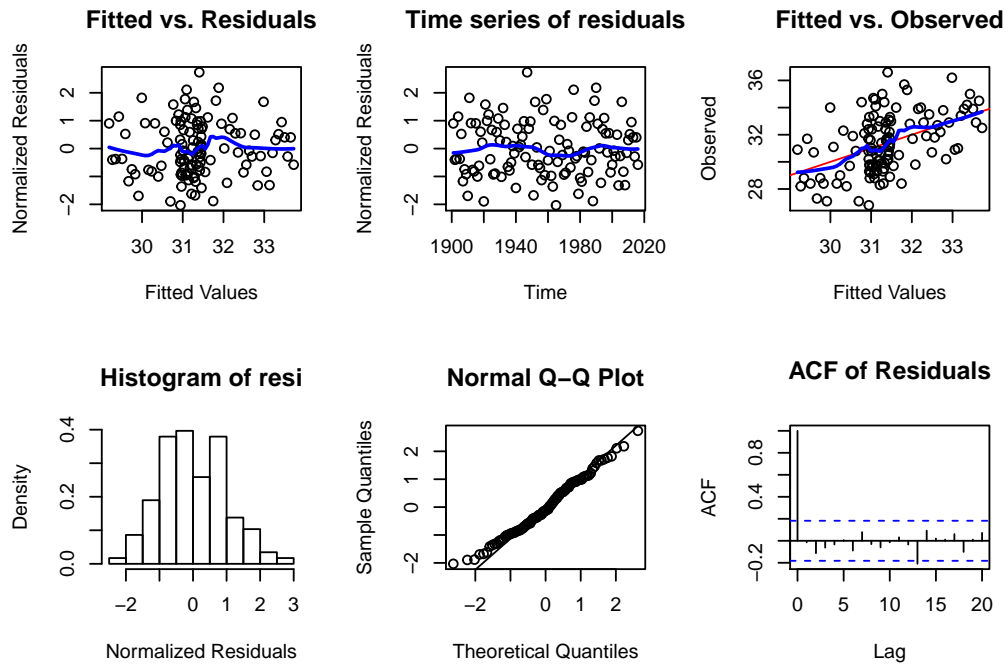


Figure C.3: Diagnostics of the chosen GAM model with MA(1) errors, based on the residuals.

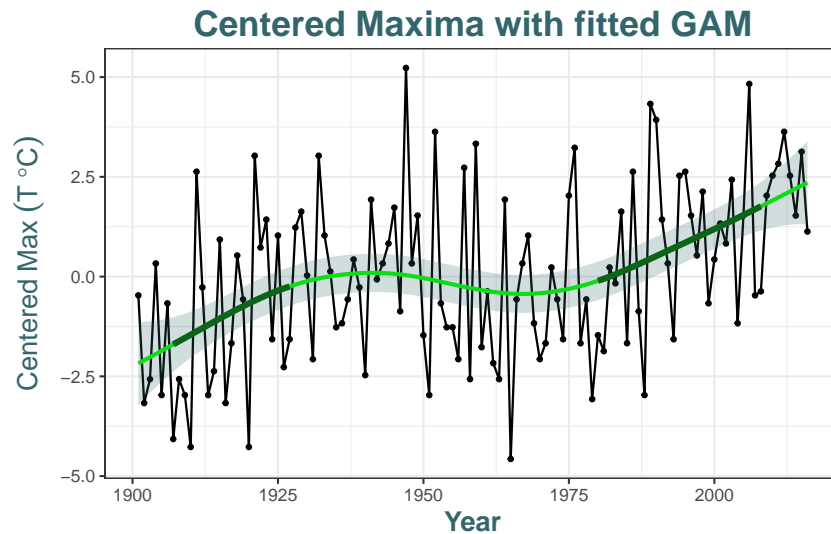


Figure C.4: Series of annual maxima together with the fitted GAM model (in green). Thicker lines indicate that the increase is significant for pointwise confidence interval. Shaded area represent the "95%" interval which looks quite narrow.

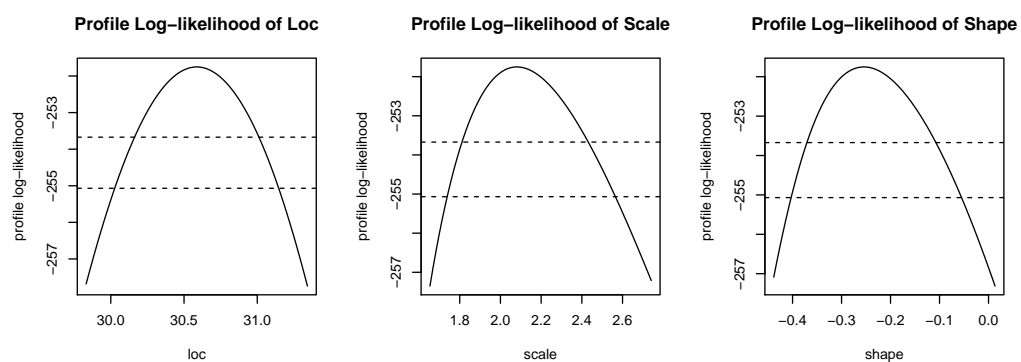


Figure C.5: The two horizontal dotted lines represent the 95% (above) and 99% (below) confidence intervals when we take the intersection on the horizontal axis.

APPENDIX D

GITHUB REPOSITORY STRUCTURE

The Github repository build for this thesis can be found on this address :

<https://github.com/proto4426/PissoortThesis>

where the R package **PissoortThesis** is located. It has the following **structure** :

- **/R/** : contains the scripts with all the functions that are made available through the package. The functions are located in the script by "category", i.e.
 - **1UsedFunc.R** : some functions created for the introduction, the stationary analysis of yearly maxima in GEV, analysis in POT, nonstationary analysis in GEV and POT,...
 - **BayesFunc.R** : functions created for the Bayesian Analysis, e.g. the Metropolis-Hastings algorithm and Gibbs Sampler (both stationary and nonstationary).
 - **BootstrapFunc.R** (not updated)
 - **NeuralNetsFunc.R** : functions slightly refined from Cannon (2010) to allow for better outputs for the nonstationary analysis with NN.
 - **runExample.R** : contains the function allowing to run the Shiny applications directly through the package (put the name of the application in ' ' in the function to load the application).

The documentation of the functions are directly made through the package, by typing `?Function_Name`.

- **/Scripts-R/**
 - **1GEV_plots(chap1).R** and **1GEV_ggplot(chap1).R** : contain the plots made for the chapter 1, but only the last latter scripts contain the code to construct the final plots (made with `ggplot2`)
 - **1intro_stationary.R** introduction and preprocessing + descriptive analysis, stationary analysis of yearly maxima in GEV, analysis in POT, analysis with other time scale and with minima, ...
 - **1intro_trends(splines).R**
 - **1intro_stationary.R**

- **1intro_stationary.R**
- **1intro_stationary.R**
- **1intro_stationary.R**
- **1intro_stationary.R**
- **1intro_stationary.R**
- **1intro_stationary.R**
- **/Shiny_app_visu/** (not updated)
- **/data/**
- **/inst/**
- **/man/**
- **/stan/**
- **/vignettes/**

Bibliography

- A. AghaKouchak, D. Easterling, K. Hsu, S. Schubert, and S. Sorooshian, editors. *Extremes in a Changing Climate*, volume 65 of *Water Science and Technology Library*. Springer Netherlands, Dordrecht, 2013. ISBN 978-94-007-4478-3 978-94-007-4479-0. URL <http://link.springer.com/10.1007/978-94-007-4479-0>.
- E. Angelino, M. J. Johnson, and R. P. Adams. Patterns of Scalable Bayesian Inference. *Foundations and Trends® in Machine Learning*, 9(2-3):119–247, 2016. ISSN 1935-8237, 1935-8245. doi: 10.1561/22000000052. URL <http://www.nowpublishers.com/article/Details/MAL-052>.
- A. A. Balkema and L. d. Haan. Residual Life Time at Great Age. *The Annals of Probability*, 2(5): 792–804, Oct. 1974. ISSN 0091-1798, 2168-894X. doi: 10.1214/aop/1176996548. URL <http://projecteuclid.org/euclid.aop/1176996548>.
- J. Beirlant, J. L. Teugels, and P. Vynckier. *Practical Analysis of Extreme Values*. Leuven University Press, 1996. ISBN 978-90-6186-768-5. Google-Books-ID: yLR3QgAACAAJ.
- J. Beirlant, Y. Goegebeur, J. Segers, and J. Teugels. *Statistics of Extremes: Theory and Applications*. John Wiley & Sons, Mar. 2006. ISBN 978-0-470-01237-6. Google-Books-ID: jqmRwfG6aloC.
- M. Betancourt. Diagnosing Suboptimal Cotangent Disintegrations in Hamiltonian Monte Carlo. *arXiv:1604.00695 [stat]*, Apr. 2016. URL <http://arxiv.org/abs/1604.00695>. arXiv: 1604.00695.
- M. Betancourt and M. Girolami. Hamiltonian Monte Carlo for hierarchical models. *Current trends in Bayesian methodology with applications*, 79:30, 2015. URL https://books.google.com/books?hl=en&lr=&id=gi6sCQAAQBAJ&oi=fnd&pg=PA79&dq=%22interactions+between+the+levels+in+the+hierarchy%22+%22those+pathologies+either+make+the+algorithms%22+%22the+pathologies+typical+of+hierarchical+models.%22+%22the+dramatic+variations+in+curvature+in+order%22+&ots=RKLEjsNycE&sig=HvrVDWcWWdoOP_JIdjz1a0On1iM.

- A. Bolívar, E. Díaz-Francés, J. Ortega, and E. Vilchis. Profile Likelihood Intervals for Quantiles in Extreme Value Distributions. *arXiv preprint arXiv:1005.3573*, 2010. URL <http://arxiv.org/abs/1005.3573>.
- A. J. Cannon. A flexible nonlinear modelling framework for nonstationary generalized extreme value analysis in hydroclimatology. *Hydrological Processes*, 24(6):673–685, Mar. 2010. ISSN 08856087. doi: 10.1002/hyp.7506. URL <http://doi.wiley.com/10.1002/hyp.7506>.
- S. Coles. *An Introduction to Statistical Modeling of Extreme Values*. Springer Series in Statistics. Springer London, London, 2001. ISBN 978-1-84996-874-4 978-1-4471-3675-0. URL <http://link.springer.com/10.1007/978-1-4471-3675-0>.
- S. G. Coles and M. J. Dixon. Likelihood-Based Inference for Extreme Value Models. *Extremes*, 2(1): 5–23, Mar. 1999. ISSN 1386-1999, 1572-915X. doi: 10.1023/A:1009905222644. URL <http://link.springer.com/article/10.1023/A:1009905222644>.
- A. C. Davison, S. A. Padoan, and M. Ribatet. Statistical Modeling of Spatial Extremes. *Statistical Science*, 27(2):161–186, 2012. ISSN 0883-4237. URL <http://www.jstor.org/stable/41714789>.
- A. L. M. Dekkers and L. D. Haan. On the Estimation of the Extreme-Value Index and Large Quantile Estimation. *The Annals of Statistics*, 17(4):1795–1832, Dec. 1989. ISSN 0090-5364, 2168-8966. doi: 10.1214/aos/1176347396. URL <http://projecteuclid.org/euclid.aos/1176347396>.
- A. L. M. Dekkers, J. H. J. Einmahl, and L. D. Haan. A Moment Estimator for the Index of an Extreme-Value Distribution. *The Annals of Statistics*, 17(4):1833–1855, Dec. 1989. ISSN 0090-5364, 2168-8966. doi: 10.1214/aos/1176347397. URL <http://projecteuclid.org/euclid.aos/1176347397>.
- J. E. Dennis and R. B. Schnabel. *Numerical Methods for Unconstrained Optimization and Nonlinear Equations*. Society for Industrial and Applied Mathematics, Philadelphia, Jan. 1987. ISBN 978-0-89871-364-0.
- D. K. Dey and J. Yan. *Extreme Value Modeling and Risk Analysis: Methods and Applications*. CRC Press, Jan. 2016. ISBN 978-1-4987-0131-0. Google-Books-ID: PYhUCwAAQBAJ.
- P. Embrechts, C. Klüppelberg, and T. Mikosch. *Modelling Extremal Events: for Insurance and Finance*. Springer Berlin Heidelberg, Feb. 2011. ISBN 978-3-642-08242-9. Google-Books-ID: dfZecgAACAAJ.
- M. Falk and F. Marohn. Von Mises Conditions Revisited. *The Annals of Probability*, 21(3):1310–1328, July 1993. ISSN 0091-1798, 2168-894X. doi: 10.1214/aop/1176989120. URL <http://projecteuclid.org/euclid.aop/1176989120>.
- L. Fawcett and D. Walshaw. Estimating return levels from serially dependent extremes: ESTIMATING RETURN LEVELS FROM SERIALY DEPENDENT EXTREMES. *Environmetrics*, 23(3):272–283, May 2012. ISSN 11804009. doi: 10.1002/env.2133. URL <http://doi.wiley.com/10.1002/env.2133>.
- R. A. Fisher and L. H. C. Tippett. Limiting Forms of the Frequency Distribution of the Largest or Smallest Member of a Sample. *ResearchGate*, 24(02):180–190, Jan. 1928. ISSN 1469-8064. doi: 10.1017/S0305004100015681. URL [https://www.researchgate.net/publication/](https://www.researchgate.net/publication/10.1017/S0305004100015681)

[230663919_Limiting_Forms_of_the_Frequency_Distribution_of_the_Largest_or_Smallest_Member_of_a_Sample.](#)

- A. Gelman, J. B. Carlin, H. S. Stern, D. B. Dunson, A. Vehtari, and D. B. Rubin. *Bayesian Data Analysis, Third Edition*. CRC Press, Nov. 2013. ISBN 978-1-4398-4095-5. Google-Books-ID: ZXKL6AQAAQBAJ.
- A. Gelman, J. Hwang, and A. Vehtari. Understanding Predictive Information Criteria for Bayesian Models. *Statistics and Computing*, 24(6):997–1016, Nov. 2014. ISSN 0960-3174. doi: 10.1007/s11222-013-9416-2. URL <http://dx.doi.org/10.1007/s11222-013-9416-2>.
- E. Gilleland and R. W. Katz. **extRemes** 2.0: An Extreme Value Analysis Package in R. *Journal of Statistical Software*, 72(8), 2016. ISSN 1548-7660. doi: 10.18637/jss.v072.i08. URL <http://www.jstatsoft.org/v72/i08/>.
- B. Gnedenko. Sur La Distribution Limite Du Terme Maximum D’Une Serie Aleatoire. *Annals of Mathematics*, 44(3):423–453, 1943. ISSN 0003-486X. doi: 10.2307/1968974. URL <http://www.jstor.org/stable/1968974>.
- J. A. Greenwood, J. M. Landwehr, N. C. Matalas, and J. R. Wallis. Probability weighted moments: Definition and relation to parameters of several distributions expressible in inverse form. *Water Resources Research*, 15(5):1049–1054, Oct. 1979. ISSN 1944-7973. doi: 10.1029/WR015i005p01049. URL <http://onlinelibrary.wiley.com/doi/10.1029/WR015i005p01049/abstract>.
- S. D. Grimshaw. Computing Maximum Likelihood Estimates for the Generalized Pareto Distribution. *Technometrics*, 35(2):185, May 1993. ISSN 00401706. doi: 10.2307/1269663. URL <http://www.jstor.org/stable/1269663?origin=crossref>.
- L. d. Haan. *On regular variation and its application to the weak convergence of sample extremes*. Mathematisch Centrum, 1970a. Google-Books-ID: sQ3vAAAAMAAJ.
- L. d. Haan and A. Ferreira. *Extreme value theory: an introduction*. Springer series in operations research. Springer, New York ; London, 2006. ISBN 978-0-387-23946-0. OCLC: ocm70173287.
- L. F. M. D. Haan. *On regular variation and its application to the weak convergence of sample extremes*. Mathematisch Centrum, 1970b. Google-Books-ID: SOKSAQAIAAJ.
- M. Hartmann and R. Ehlers. Bayesian Inference for Generalized Extreme Value Distributions via Hamiltonian Monte Carlo. *Communications in Statistics - Simulation and Computation*, pages 0–0, Mar. 2016. ISSN 0361-0918, 1532-4141. doi: 10.1080/03610918.2016.1152365. URL <http://arxiv.org/abs/1410.4534>. arXiv: 1410.4534.
- T. Hastie and R. Tibshirani. Generalized Additive Models. *Statistical Science*, 1(3):297–310, 1986. ISSN 0883-4237. URL <http://www.jstor.org/stable/2245459>.
- W. K. Hastings. Monte Carlo sampling methods using Markov chains and their applications. *Biometrika*, 57(1):97–109, Apr. 1970. ISSN 0006-3444. URL <https://academic.oup.com/biomet/article-abstract/57/1/97/2721936/Monte-Carlo-sampling-methods-using-Markov-chains>.

- J. R. M. Hosking and J. R. Wallis. Parameter and Quantile Estimation for the Generalized Pareto Distribution. *Technometrics*, 29(3):339, Aug. 1987. ISSN 00401706. doi: 10.2307/1269343. URL <http://www.jstor.org/stable/1269343?origin=crossref>.
- J. R. M. J. R. M. Hosking and J. R. Wallis. *Regional frequency analysis : an approach based on L-moments*. Cambridge ; New York : Cambridge University Press, 1997. ISBN 0521430453 (hardbound).
- J. P. Iii. Statistical Inference Using Extreme Order Statistics. *The Annals of Statistics*, 3(1):119–131, Jan. 1975. ISSN 0090-5364, 2168-8966. doi: 10.1214/aos/1176343003. URL <http://projecteuclid.org/euclid.aos/1176343003>.
- A. M. G. Klein Tank, J. B. Wijngaard, G. P. Können, R. Böhm, G. Demarée, A. Gocheva, M. Mileta, S. Pashiardis, L. Hejkrlik, C. Kern-Hansen, R. Heino, P. Bessemoulin, G. Müller-Westermeier, M. Tzanakou, S. Szalai, T. Pálsdóttir, D. Fitzgerald, S. Rubin, M. Capaldo, M. Maugeri, A. Leitass, A. Bukantis, R. Aberfeld, A. F. V. van Engelen, E. Forland, M. Miletus, F. Coelho, C. Mares, V. Razuvaev, E. Nieplova, T. Cegnar, J. Antonio López, B. Dahlström, A. Moberg, W. Kirchhofer, A. Ceylan, O. Pachaliuk, L. V. Alexander, and P. Petrovic. Daily dataset of 20th-century surface air temperature and precipitation series for the European Climate Assessment. *International Journal of Climatology*, 22(12):1441–1453, Oct. 2002. ISSN 1097-0088. doi: 10.1002/joc.773. URL <http://onlinelibrary.wiley.com/doi/10.1002/joc.773/abstract>.
- A. N. Kolmogorov, N. Morrison, and A. T. Bharucha-Reid. *Foundations of the theory of probability*. Chelsea Publishing Company, New York, 1956. OCLC: 751236060.
- M. R. Leadbetter. On extreme values in stationary sequences. *Zeitschrift für Wahrscheinlichkeitstheorie und Verwandte Gebiete*, 28(4):289–303, Dec. 1974. ISSN 0044-3719, 1432-2064. doi: 10.1007/BF00532947. URL <http://link.springer.com/article/10.1007/BF00532947>.
- M. R. Leadbetter, G. Lindgren, and H. Rootzén. *Extremes and Related Properties of Random Sequences and Processes*. Springer Series in Statistics. Springer New York, New York, NY, 1983. ISBN 978-1-4612-5451-5 978-1-4612-5449-2. URL <http://link.springer.com/10.1007/978-1-4612-5449-2>.
- F. Leisch. Creating r packages: A tutorial. 2008. URL <https://epub.ub.uni-muenchen.de/6175/>.
- A. A. Lindsey and J. E. Newman. Use of Official Wather Data in Spring Time: Temperature Analysis of an Indiana Phenological Record. *Ecology*, 37(4):812–823, Oct. 1956. ISSN 1939-9170. doi: 10.2307/1933072. URL <http://onlinelibrary.wiley.com/doi/10.2307/1933072/abstract>.
- Y. Liu, A. Gelman, and T. Zheng. Simulation-efficient shortest probability intervals. *Statistics and Computing*, 25(4):809–819, July 2015. ISSN 0960-3174, 1573-1375. doi: 10.1007/s11222-015-9563-8. URL <http://link.springer.com/10.1007/s11222-015-9563-8>.
- J. Maindonald and J. Braun. *Data analysis and graphics using R: an example-based approach*, volume 10. Cambridge University Press, 2006. URL <https://books.google.com/books?hl=en&lr=&id=d70eVD6SKBsC&oi=fnd&pg=PA5&dq=%22methods,+stats,+graphics,%22+%22+.+.+.%22+%22+.+.+.%22+%22+.+.+.>

+.+. %22+%22.+.+. %22+%22.+.+.+.+. %22+&ots=2mOJGySmmE&sig=MLMmYfuz4lg9QmCsb2c6DnE6ShE.

- G. Marra and S. N. Wood. Coverage properties of confidence intervals for generalized additive model components. *Scandinavian Journal of Statistics*, 39(1):53–74, 2012. URL <http://onlinelibrary.wiley.com/doi/10.1111/j.1467-9469.2011.00760.x/full>.
- P. C. D. Milly, J. Betancourt, M. Falkenmark, R. M. Hirsch, Z. W. Kundzewicz, D. P. Lettenmaier, and R. J. Stouffer. Climate change. Stationarity is dead: whither water management? *Science (New York, N.Y.)*, 319(5863):573–574, Feb. 2008a. ISSN 1095-9203. doi: 10.1126/science.1151915.
- P. C. D. Milly, J. Betancourt, M. Falkenmark, R. M. Hirsch, Z. W. Kundzewicz, D. P. Lettenmaier, and R. J. Stouffer. Stationarity Is Dead: Whither Water Management? *Science*, 319(5863):573–574, Feb. 2008b. ISSN 0036-8075, 1095-9203. doi: 10.1126/science.1151915. URL <http://science.sciencemag.org/content/319/5863/573>.
- M. Mudelsee. *Climate Time Series Analysis*, volume 51 of *Atmospheric and Oceanographic Sciences Library*. Springer International Publishing, Cham, 2014. ISBN 978-3-319-04449-1 978-3-319-04450-7. URL <http://link.springer.com/10.1007/978-3-319-04450-7>.
- R. M. Neal and others. MCMC using Hamiltonian dynamics. *Handbook of Markov Chain Monte Carlo*, 2:113–162, 2011. URL [https://books.google.com/books?hl=en&lr=&id=qfRsAIKZ4rIC&oi=fnd&pg=PA113&dq=%22in+the+following+decades+\(see+Frenkel+and+Smit,%22+%22implementation.+The+%E2%80%9C9Cleafrog%E2%80%9D+scheme+that+is+typically+used+is+quite%22+%22implemented+with+the+leapfrog+method.+A+state+proposed+in+this+way+can%22+&ots=RbAabV3a4S&sig=WFFHl6fmkouIgzGmizb8pnxBPVM](https://books.google.com/books?hl=en&lr=&id=qfRsAIKZ4rIC&oi=fnd&pg=PA113&dq=%22in+the+following+decades+(see+Frenkel+and+Smit,%22+%22implementation.+The+%E2%80%9C9Cleafrog%E2%80%9D+scheme+that+is+typically+used+is+quite%22+%22implemented+with+the+leapfrog+method.+A+state+proposed+in+this+way+can%22+&ots=RbAabV3a4S&sig=WFFHl6fmkouIgzGmizb8pnxBPVM).
- S. Ni and D. Sun. Noninformative priors and frequentist risks of bayesian estimators of vector-autoregressive models. *Journal of Econometrics*, 115(1):159–197, July 2003. ISSN 0304-4076. doi: 10.1016/S0304-4076(03)00099-X. URL <http://www.sciencedirect.com/science/article/pii/S030440760300099X>.
- P. J. Northrop and N. Attalides. Posterior propriety in Bayesian extreme value analyses using reference priors. *Statistica Sinica*, 26(2), Apr. 2016. ISSN 1017-0405. URL <http://dx.doi.org/10.5705/ss.2014.034>.
- M. C. Peel, Q. J. Wang, R. M. Vogel, and T. A. McMAHON. The utility of L-moment ratio diagrams for selecting a regional probability distribution. *Hydrological Sciences Journal*, 46(1):147–155, 2001. URL <http://www.tandfonline.com/doi/abs/10.1080/02626660109492806>.
- E. C. Pinheiro and S. L. P. Ferrari. A comparative review of generalizations of the Gumbel extreme value distribution with an application to wind speed data. *arXiv:1502.02708 [stat]*, Feb. 2015. URL <http://arxiv.org/abs/1502.02708>. arXiv: 1502.02708.
- R.-D. Reiss and M. Thomas. *Statistical analysis of extreme values: with applications to insurance, finance, hydrology and other fields ; [includes CD-ROM]*. Birkhäuser, Basel, 3. ed edition, 2007. ISBN 978-3-7643-7230-9 978-3-7643-7399-3. OCLC: 180885018.

- S. I. Resnick. *Extreme Values, Regular Variation and Point Processes*. Springer Series in Operations Research and Financial Engineering. Springer New York, New York, NY, 1987. ISBN 978-0-387-75952-4 978-0-387-75953-1. URL <http://link.springer.com/10.1007/978-0-387-75953-1>.
- M. Ribatet. A User's Guide to the POT Package (Version 1.4). *month*, 2006. URL <http://www.unalmed.edu.co/~ndgiraldo/Archivos%20Lectura/Archivos%20curso%20Riesgo%20Operativo/POT.pdf>.
- M. Ribatet, C. Dombry, and M. Oesting. Spatial Extremes and Max-Stable Processes. In *Extreme Value Modeling and Risk Analysis: Methods and Applications*, pages 179–194. Chapman and Hall/CRC, 2015. URL https://books.google.com/books?hl=en&lr=&id=PYhUCwAAQBAJ&oi=fnd&pg=PA179&dq=%22rainfall+amount+in+this+catchment,+i.e.,+evaluating+probabilities%22+%22Extremes+and+max-stable%22+%22as+far+as+spatial+extremes+are+of+concern.+Similarly+to+the+univariate%22+%22To+be+consistent+with+the+univariate+extreme+value+theory,+Theorem%22+&ots=phZY14rcy5&sig=-_n3q707dYkioxyG-2QpsYB62jk.
- P. Ribereau, E. Masiello, and P. Naveau. Skew generalized extreme value distribution: Probability-weighted moments estimation and application to block maxima procedure. *Communications in Statistics - Theory and Methods*, 45(17):5037–5052, Sept. 2016. ISSN 0361-0926. doi: 10.1080/03610926.2014.935434. URL <http://dx.doi.org/10.1080/03610926.2014.935434>.
- G. O. Roberts, A. Gelman, and W. R. Gilks. Weak convergence and optimal scaling of random walk Metropolis algorithms. *The Annals of Applied Probability*, 7(1):110–120, Feb. 1997. ISSN 1050-5164, 2168-8737. doi: 10.1214/aoap/1034625254. URL <http://projecteuclid.org/euclid.aoap/1034625254>.
- G. Rosso. Extreme Value Theory for Time Series using Peak-Over-Threshold method. *arXiv preprint arXiv:1509.01051*, 2015. URL <http://arxiv.org/abs/1509.01051>.
- D. Ruppert, M. P. Wand, and R. J. Carroll. *Semiparametric Regression*. Cambridge Books, Cambridge University Press, 2003. URL <http://econpapers.repec.org/bookchap/cupcbooks/9780521785167.htm>.
- J. Segers. Generalized Pickands Estimators for the Extreme Value Index: Minimal Asymptotic Variance and Bias Reduction. 2001. URL <http://citeseerx.ist.psu.edu/viewdoc/citations;jsessionid=0C9C7CAE3938CA7A9B280DEA739EFDA9?doi=10.1.1.7.1713>.
- B. A. Shaby, B. J. Reich, D. Cooley, and C. G. Kaufman. A Markov-switching model for heat waves. *The Annals of Applied Statistics*, 10(1):74–93, Mar. 2016. ISSN 1932-6157. doi: 10.1214/15-AOAS873. URL <http://projecteuclid.org/euclid.aoas/1458909908>.
- C. Sherlock, G. Roberts, and others. Optimal scaling of the random walk Metropolis on elliptically symmetric unimodal targets. *Bernoulli*, 15(3):774–798, 2009. URL <http://projecteuclid.org/euclid.bj/1251463281>.

- V. P. Singh and H. Guo. Parameter estimation for 3-parameter generalized pareto distribution by the principle of maximum entropy (POME). *Hydrological Sciences Journal*, 40(2):165–181, Apr. 1995. ISSN 0262-6667, 2150-3435. doi: 10.1080/02626669509491402. URL <http://www.tandfonline.com/doi/abs/10.1080/02626669509491402>.
- R. L. Smith. Maximum likelihood estimation in a class of nonregular cases. *Biometrika*, 72(1):67–90, Apr. 1985. ISSN 0006-3444. doi: 10.1093/biomet/72.1.67. URL <https://academic.oup.com/biomet/article-abstract/72/1/67/242523/Maximum-likelihood-estimation-in-a-class-of>.
- D. J. Spiegelhalter, N. G. Best, B. P. Carlin, and A. Van Der Linde. Bayesian measures of model complexity and fit. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 64(4):583–639, 2002. URL <http://onlinelibrary.wiley.com/doi/10.1111/1467-9868.00353/full>.
- M.-S. Suh, S.-G. Oh, D.-K. Lee, D.-H. Cha, S.-J. Choi, C.-S. Jin, and S.-Y. Hong. Development of New Ensemble Methods Based on the Performance Skills of Regional Climate Models over South Korea. *Journal of Climate*, 25(20):7067–7082, May 2012. ISSN 0894-8755. doi: 10.1175/JCLI-D-11-00457.1. URL <http://journals.ametsoc.org/doi/abs/10.1175/JCLI-D-11-00457.1>.
- A. Vehtari, A. Gelman, and J. Gabry. Practical Bayesian model evaluation using leave-one-out cross-validation and WAIC. *Statistics and Computing*, Aug. 2016. ISSN 0960-3174, 1573-1375. doi: 10.1007/s11222-016-9696-4. URL <http://arxiv.org/abs/1507.04544>. arXiv: 1507.04544.
- R. Von Mises. La distribution de la plus grande de n valeurs. *Rev., Math, Union Interbalcanique*, 1: pp.141–160, 1936.
- R. Wada, T. Waseda, and P. Jonathan. Extreme value estimation using the likelihood-weighted method. *Ocean Engineering*, 124:241–251, 2016. ISSN 0029-8018. doi: 10.1016/j.oceaneng.2016.07.063.
- J. L. Wadsworth. Exploiting Structure of Maximum Likelihood Estimators for Extreme Value Threshold Selection. *Technometrics*, 58(1):116–126, Jan. 2016. ISSN 0040-1706. doi: 10.1080/00401706.2014.998345. URL <http://dx.doi.org/10.1080/00401706.2014.998345>.
- Q. J. Wang. LH moments for statistical analysis of extreme events. *Water Resources Research*, 33(12):2841–2848, Dec. 1997. ISSN 1944-7973. doi: 10.1029/97WR02134. URL <http://onlinelibrary.wiley.com/doi/10.1029/97WR02134/abstract>.
- S. Watanabe. Asymptotic equivalence of Bayes cross validation and widely applicable information criterion in singular learning theory. *Journal of Machine Learning Research*, 11(Dec):3571–3594, 2010. URL <http://www.jmlr.org/papers/v11/watanabe10a.html>.
- R. Yang and J. O. Berger. *A catalog of noninformative priors*. Institute of Statistics and Decision Sciences, Duke University, 1996. URL <http://www.stats.org.uk/priors/noninformative/YangBerger1998.pdf>.
- C. Zhou. The extent of the maximum likelihood estimator for the extreme value index. *Journal of Multivariate Analysis*, 101(4):971–983, Apr. 2010. ISSN 0047-259X. doi: 10.1016/j.

jmva.2009.09.013. URL <http://www.sciencedirect.com/science/article/pii/S0047259X0900178X>.