



UNIVERSITÀ  
DEGLI STUDI  
FIRENZE

SCUOLA DI ECONOMIA E MANAGEMENT

LAUREA MAGISTRALE IN STATISTICA, SCIENZE ATTUARIALI E  
FINANZIARIE

**Extreme value theory: statistical  
estimates and actuarial evaluations for  
seismic risk in Italy**

Anno Accademico 2015/2016

*Relatore*

Prof. Luigi Vannucci

*Candidato*

Francesca Giorgolo

*Correlatore*

Prof. Marcello Galeotti



<b>Introduction</b>	<b>1</b>
<b>1 Classical extreme value theory</b>	<b>4</b>
1.1 Basic definitions . . . . .	4
1.2 Limit law for maxima . . . . .	5
1.3 Maximum domain of attraction . . . . .	11
1.3.1 Case of Fréchet distribution $\Phi_\alpha$ . . . . .	13
1.3.2 Case of Weibull distribution $\Psi_\alpha$ . . . . .	14
1.3.3 Case of Gumbel distribution $\Lambda$ . . . . .	15
1.4 Generalised extreme value distribution . . . . .	16
1.5 Inference for GEV distribution . . . . .	20
1.6 Threshold models . . . . .	24
1.7 Inference for threshold models . . . . .	28
<b>2 Statistical analysis of Italian earthquakes data</b>	<b>30</b>
2.1 Introduction . . . . .	30
2.2 Preparation of the dataset . . . . .	32
2.3 GEV distribution parameters estimation . . . . .	34
2.3.1 Global analysis . . . . .	34
2.3.2 Ten years blocks . . . . .	37
2.3.3 Fifty years blocks . . . . .	39
2.3.4 Zone analysis . . . . .	41
2.3.5 Results comparison . . . . .	45
2.4 GPD parameters estimation . . . . .	46
2.4.1 Global analysis . . . . .	46
2.4.2 Zone analysis . . . . .	50

2.4.3	Results comparison . . . . .	63
<b>3</b>	<b>Actuarial evaluations</b>	<b>64</b>
3.1	Introduction . . . . .	64
3.2	Recent catastrophic earthquakes . . . . .	64
3.3	Insurer's loss due to catastrophic events . . . . .	66
3.4	Premium and equilibrium reserve . . . . .	72
3.5	Premium diversification among zones . . . . .	75
	<b>Conclusions</b>	<b>79</b>

Man has always felt the need to protect himself from events which could affect the peaceful flow of his life and a solution he found was transferring risk to someone else less averse to risk than him: the insurer.

However insurance companies are not willing to insure every type of risk because they have to protect themselves from failure; in particular natural disasters pose several challenges to insurers because they involve potentially high losses that are extremely uncertain.

In this thesis we will focus on a specific topic: the seismic risk.

In the FIRST CHAPTER we will illustrate the classical extreme value theory. Considering a sequence  $X_1, X_2, X_3, \dots$  of independent identically distributed non-degenerate random variables with common distribution function  $F$  we will define the sample maxima  $M_n$  as the maximum of  $X_1, \dots, X_n$ ,  $n \geq 1$ , and we will prove that the limit distribution of normalised maxima belongs to one of three distribution function families: the Gumbel  $\Lambda$ , the Fréchet  $\Phi_\alpha$  and the Weibull  $\Psi_\alpha$  with  $\alpha > 0$ , called standard extreme value distributions.

After we will introduce the concept of maximum domain of attraction of the extreme value distribution, characterizing it on varying of the distribution family: in particular will be the tail of the distribution function  $F$  to make the difference.

From a statistical point of view is simpler dealing with a single distribution function for maxima instead of three different distribution families, thus we will introduce the generalised extreme value distribution (GEV)  $H_{\xi;\mu,\sigma}$  where  $\mu \in \mathbb{R}$  is the location parameter,  $\sigma > 0$  is the scale parameter and  $\xi \in \mathbb{R}$  is the shape parameter which value, depending whether positive, negative or

zero, distinguishes the three families. Then we will show how making inference for this distribution, in order to estimate its parameters and check the model adequacy.

After that we will introduce the generalised Pareto distribution (GPD)  $G_{\xi, \tilde{\mu}, \tilde{\sigma}}$ , which is the limit distribution of the scaled excess over an high threshold  $u$ , that is  $((X - u)/a(u)|X > u)$ . This distribution has some interesting properties, in particular for  $\tilde{\mu} = 0$  it well approximates the excess distribution function  $F_u(x) = \mathbb{P}(X - u \leq x|X > u)$ ,  $x \geq 0$ . Finally we will show how to make inference also for parameters of this distribution.

In the SECOND CHAPTER we will apply the theory developed in the first chapter to Italian earthquakes data available in the Parametric Catalogue of Italian Earthquakes (CPTI15). Our purpose is to find the distribution family of the maximum magnitude.

First of all we will clean the dataset. Since records refer to the period 1000-2014, we will keep only reliable and satisfactory records; instead records referring to replicates will be dropped, since results obtained in the first chapter concern only independent events.

Then we will estimate GEV distribution parameters for the whole Italian territory, first for annual block maxima, then for ten years block maxima and finally for fifty years block maxima, checking also model adequacy. After that we will divide the Italian territory in six regions, *Sea/Foreign*, *Alps*, *Po valley*, *Centre*, *South* and *Islands*, distinguished by their seismic, and we will analyse what changes estimating zone by zone GEV distribution parameters. Finally we will use the GPD approach to investigate values assumed by the shape parameter  $\xi$ , first choosing a threshold value for all Italian records and then for regional records, every time checking the adequacy of the fitted model.

In the THIRD CHAPTER we will use estimated distribution functions to set premiums for seismic risk coverage. Some assumptions will be needed, in particular we will assume that an insurer insures the whole Italian territory, consisting of 27 million housing units with an overall reconstruction cost of 3900 billion Euro. Furthermore we will assume that the loss depends in a proportional manner on the energy released by the earthquake, that in turn depends on the annual maximum magnitude value. So in our context the loss  $L$  will be a random variable depending on value of annual block maxima  $M_w$ .

Then, from costs of seven recent catastrophic earthquakes concentrated in

the period 1968-2014, we will estimate the proportionality constant  $k$  which will let us infer the extent of loss for all annual block maxima.

Using quantiles of the  $M_w$ 's distribution function, we will be able to relate the loss with its exceedance probability and return period: in particular we will calculate the expected annual loss and the premium, first for the whole Italian territory then differentiating by region.

# CHAPTER 1

## CLASSICAL EXTREME VALUE THEORY

### 1.1 Basic definitions

Let  $(\Omega, \mathfrak{S}, \mathbb{P})$  a probability space, where  $\Omega$  is the events set,  $\mathfrak{S}$  is a  $\sigma$ -algebra on  $\Omega$  and  $\mathbb{P} : \mathfrak{S} \rightarrow [0, 1]$  is a probability measure.

**Definition 1.1.1.** Let  $X : \Omega \rightarrow \mathbb{R}$  be a random variable. The function  $F$  defined by

$$\begin{aligned} F : \mathbb{R} &\rightarrow [0, 1] \\ x &\mapsto \mathbb{P}(] - \infty, x]) = \mathbb{P}(X \leq x). \end{aligned}$$

is called *distribution function*.

The *tail* of the distribution function is  $\bar{F} = 1 - F$ .

**Definition 1.1.2.** Suppose  $h$  is a non-decreasing function on  $\mathbb{R}$ . The function defined as

$$h^{\leftarrow}(t) = \inf \{x \in \mathbb{R} : h(x) \geq t\}$$

is called *generalised inverse* of  $h$ , with the convention that  $\inf \{\emptyset\} = \infty$ .

**Definition 1.1.3.** The generalised inverse of the distribution function  $F$

$$F^{\leftarrow}(t) = \inf \{x \in \mathbb{R} : F(x) \geq t\}, \quad 0 < t < 1$$

is called *quantile function*.

The quantity  $x_t = F^{\leftarrow}(t)$  defines the *t-quantile* of  $F$ .



We denote by  $x_F = \sup \{x \in \mathbb{R} : F(x) \leq 1\}$  the *right endpoint* of  $F$ .

It can be useful also remember different kind of convergence.

**Definition 1.1.4.** We say a sequence of real random variables  $\{X_n\}_{n \geq 1}$  converges to a random variable  $X$  on  $\mathbb{R}$

- **in probability** if  $\forall \epsilon > 0 \quad \lim_{n \rightarrow \infty} \mathbb{P}(|X_n - X| > \epsilon) = 0$ ;
- **almost sure** if, for almost any  $\omega$ , the sequence of numbers  $X_n(\omega)$  converges to  $X(\omega)$ ;
- **in  $L_2$**  if every  $X_n$  (and  $X$ ) belongs to  $L_2$  and  $\lim_{n \rightarrow \infty} \|X_n - X\|_2 = \lim_{n \rightarrow \infty} \mathbb{E}[|X_n - X|^2] = 0$ .

Between them these relations hold:

- if  $X_n \xrightarrow{L_2} X$ , then  $X_n \xrightarrow{\mathbb{P}} X$ ;
- if  $X_n \xrightarrow{\text{a.s.}} X$ , then  $X_n \xrightarrow{\mathbb{P}} X$ ;
- if  $X_n \xrightarrow{\mathbb{P}} X$ , then exists a subsequence  $X_{n_k}$  such that  $X_{n_k} \xrightarrow{\text{a.s.}} X$ .

**Definition 1.1.5.** A sequence of real random variables  $\{X_n\}_{n \geq 1}$  is said to *converge in distribution* to a random variable  $X$  on  $\mathbb{R}$  if for all  $f$  continuous and bounded function  $\lim_{n \rightarrow \infty} \mathbb{E}[f(X_n)] = \mathbb{E}[f(X)]$ .

## 1.2 Limit law for maxima

Let  $X_1, X_2, X_3, \dots$  be a sequence of independent identically distributed non-degenerate random variables with common distribution function  $F$ .

**Definition 1.2.1.** We define the *sample maxima* as

$$M_n = \max\{X_1, \dots, X_n\}, \quad n \geq 1$$

and the *sample minima* as

$$\min\{X_1, \dots, X_n\} = -\max\{-X_1, \dots, -X_n\}, \quad n \geq 1.$$

Thanks to the independence and equal distribution of variables, we can easily write down the exact distribution function of  $M_n$ :

$$\begin{aligned}\mathbb{P}(M_n \leq x) &= \mathbb{P}(X_1 \leq x, \dots, X_n \leq x) \\ &= \mathbb{P}(X_1 \leq x) \cdots \mathbb{P}(X_n \leq x) \\ &= F(x) \cdots F(x) \\ &= F^n(x), \quad x \in \mathbb{R}, n \in \mathbb{N}.\end{aligned}$$

*Remark 1.2.1.* In practice this is not so useful because the common distribution function  $F$  is unknown.

We can avoid this problem using standard classical statistical techniques to estimate  $F$  from observed data, but we must keep in mind that small discrepancies in the estimate could lead to substantial discrepancies for  $F^n$ .

Alternatively we can directly look for an approximate of  $F^n$  based on extreme data only.

From this result we obtain that

- $\forall x < x_F \lim_{n \rightarrow \infty} \mathbb{P}(M_n \leq x) = \lim_{n \rightarrow \infty} F^n(x) = 0$ ;
- $\forall x \geq x_F$ , obviously in the case  $x_F < \infty$ ,  $\mathbb{P}(M_n \leq x) = F^n(x) = 1$ ;

thus

$$M_n \xrightarrow{\mathbb{P}} x_F \quad \text{where} \quad x_F \leq \infty.$$

Moreover, since the sequence of numbers  $M_n$  is non decreasing in  $n$ , we have

$$M_n \xrightarrow{\text{a.s.}} x_F$$

but this isn't enough.

In fact, our aim is showing the limit distribution of maxima but, to achieve this purpose, we need first give conditions on  $F$  under which

$$\lim_{n \rightarrow \infty} \mathbb{P}(M_n \leq u_n)$$

exists for an appropriate constant  $u_n$ .

The issue is, in contrast with sums for which the Central Limit Theorem ensures the convergence to the Normal distribution under the general condition  $\mathbb{E}[X^2] < \infty$ , that in the case of maxima we always need conditions on the tail  $\bar{F}$ ; these conditions are given in the following

**Proposition 1.2.1** (Poisson approximation).

Given a sequence  $\{u_n\} \in \mathbb{R}$  and  $\tau \in [0, \infty]$ , we have

$$\lim_{n \rightarrow \infty} n\bar{F}(u_n) \rightarrow \tau \quad \Leftrightarrow \quad \lim_{n \rightarrow \infty} \mathbb{P}(M_n \leq u_n) \rightarrow e^{-\tau}.$$

*Proof.* ( $\Rightarrow$ ) If the first equation holds then

$$\mathbb{P}(M_n \leq u_n) = F^n(u_n) = (1 - \bar{F}(u_n))^n = \left(1 - \frac{\tau}{n} + o\left(\frac{1}{n}\right)\right)^n \xrightarrow{n \rightarrow \infty} e^{-\tau}.$$

( $\Leftarrow$ ) If  $\mathbb{P}(M_n \leq u_n) \rightarrow e^{-\tau}$  then  $\bar{F}(u_n) \rightarrow 0$  and taking logarithms of this condition we have

$$\begin{aligned} \ln \mathbb{P}(M_n \leq u_n) &\rightarrow \ln(e^{-\tau}) \\ \ln(1 - \bar{F}(u_n))^n &\rightarrow -\tau \\ -n \ln(1 - \bar{F}(u_n)) &\rightarrow \tau. \end{aligned}$$

Since  $-\ln(1 - x) \sim x$  for  $x \rightarrow 0$ , then  $n\bar{F}(u_n) = \tau + o(1)$ .

More detailed proof, included the case  $\tau = \infty$ , can be found in Embrechts et al. [1], Proposition 3.1.1.  $\square$

*Remark 1.2.2.* Assume  $\tau \in (0, \infty)$  and consider the random variable

$$B_n = \sum_{i=1}^n I_{\{X_i > u_n\}}$$

which represents the number of excesses over the threshold  $u_n$ .

The indicators

$$I_{\{X_i > u_n\}} = \begin{cases} 1 & \text{if } X_i > u_n \\ 0 & \text{if } X_i \leq u_n \end{cases}, \quad i = 1, \dots, n$$

are independent Bernoulli variables with parameter  $\bar{F}(u_n)$ , then  $B_n$  is a binomial random variable such that  $B_n \sim \mathcal{B}(n, \bar{F}(u_n))$ .

In the context of extremal events  $n \rightarrow \infty$  and  $\bar{F}(u_n) \rightarrow 0$  so, applying the Poisson limit theorem we have

$$\mathbb{E}[B_n] = n\bar{F}(u_n) \rightarrow \tau \quad \Leftrightarrow \quad B_n \xrightarrow{d} \mathcal{P}(\tau).$$

In particular  $\mathbb{P}(M_n \leq u_n) = \mathbb{P}(B_n = 0) \rightarrow e^{-\tau}$ .

**Definition 1.2.2.** The distribution of a non-degenerate random variable  $X$  is called *max-stable distribution* if it satisfies

$$M_n = \max\{X_1, \dots, X_n\} \stackrel{d}{=} c_n X + d_n \quad (1.1)$$

for iid random variables  $X, X_1, \dots, X_n$  with  $n \geq 2$  and for appropriate *centring and normalising constants*  $c_n > 0, d_n \in \mathbb{R}$  respectively.

Rewriting (1.1) as

$$M_n^* := \frac{M_n - d_n}{c_n} \stackrel{d}{=} X,$$

where  $M_n^*$  is called *normalised maxima*, we can conclude that every max-stable distribution is a limit distribution for maxima of iid random variables. Furthermore, max-stable distributions are the only limit laws for  $M_n^*$ ; in fact the following theorem holds:

**Theorem 1.2.2.**

*The class of all possible non degenerate limit laws for normalised maxima  $M_n^*$  of iid random variables coincides with the class of max-stable distributions.*

*Proof.* We already proved that every max-stable distribution is a limit distribution for maxima of iid random variables.

Conversely, to prove that the limit distribution is max-stable, we assume that

$$\lim_{n \rightarrow \infty} F^n(c_n x + d_n) = H(x), \quad x \in \mathbb{R}$$

for some non-degenerate distribution function  $H$ .

Then for every  $k \in \mathbb{N}$

$$\lim_{n \rightarrow \infty} F^{nk}(c_n x + d_n) = \left( \lim_{n \rightarrow \infty} F^n(c_n x + d_n) \right)^k = H^k(x), \quad x \in \mathbb{R}.$$

Furthermore

$$\lim_{n \rightarrow \infty} F^{nk}(c_{nk} x + d_{nk}) = H(x), \quad x \in \mathbb{R}.$$

For Convergence to types theorem (see Resnick [2], Proposition 0.2) there exist  $\tilde{c}_k > 0, \tilde{d}_k \in \mathbb{R}$  such that

$$\lim_{n \rightarrow \infty} \frac{c_{nk}}{c_n} = \tilde{c}_k, \quad \lim_{n \rightarrow \infty} \frac{d_{nk} - d_n}{c_n} = \tilde{d}_k, \quad t > 0.$$

Thus, for iid random variables  $Y_1, \dots, Y_k$  with distribution function  $H$ , we obtain

$$\max\{Y_1, \dots, Y_k\} \stackrel{d}{=} \tilde{c}_k Y_1 + \tilde{d}_k.$$

□

The following theorem is the basis of classical extreme value theory, because it gives the entire range of possible limit distributions for  $M_n^*$ .

**Theorem 1.2.3** (Fisher-Tippett (1928)).

Let  $\{X_n\}$  be a sequence of non degenerate iid random variables.

If there exist norming constants  $c_n > 0$ ,  $d_n \in \mathbb{R}$  and some non degenerate distribution function  $H$  such that

$$M_n^* := \frac{M_n - d_n}{c_n} \xrightarrow{d} H,$$

then  $H$  belongs to one of the following distribution functions families:

1. **Gumbel**  $\Lambda(x) = \exp(-e^{-x}) \quad x \in \mathbb{R}$
2. **Fréchet**  $\Phi_\alpha(x) = \begin{cases} 0 & \text{if } x \leq 0 \\ \exp(-x^{-\alpha}) & \text{if } x > 0 \end{cases} \quad \alpha > 0$
3. **Weibull**  $\Psi_\alpha(x) = \begin{cases} \exp(-(-x)^\alpha) & \text{if } x \leq 0 \\ 1 & \text{if } x > 0 \end{cases} \quad \alpha > 0.$

These distribution functions are called standard extreme value distributions.

*Proof.* The proof is technical and we only sketch it.

For any  $t > 0$  we have

$$F^{[nt]}(c_{[nt]}x + d_{[nt]}) \rightarrow H(x), \quad x \in \mathbb{R},$$

where  $[t]$  denotes the integer part of  $t$ , and also

$$F^{[nt]}(c_n x + d_n) = (F^n(c_n x + d_n))^{[nt]/n} \rightarrow H^t(x), \quad x \in \mathbb{R}.$$

So, from Convergence to types theorem (see Resnick [2], Proposition 0.2 ), there exists functions  $\gamma(t) > 0$ ,  $\delta(t) \in \mathbb{R}$  such that

$$\lim_{n \rightarrow \infty} \frac{c_n}{c_{[nt]}} = \gamma(t), \quad \lim_{n \rightarrow \infty} \frac{d_n - d_{[nt]}}{c_{[nt]}} = \delta(t), \quad t > 0,$$

and

$$H^t(x) = H(\gamma(t)x + \delta(t)).$$

On one hand, for  $s > 0$ , we have

$$H^{st}(x) = H(\gamma(st)x + \delta(st))$$

and on the other hand

$$\begin{aligned}
H^{st}(x) &= (H^s(x))^t \\
&= H^t(\gamma(s)x + \delta(s)) \\
&= H(\gamma(t)(\gamma(s)x + \delta(s)) + \delta(t)) \\
&= H(\gamma(t)\gamma(s)x + \gamma(t)\delta(s) + \delta(t)).
\end{aligned}$$

Since  $G$  is assumed non-degenerate we therefore conclude for  $t > 0$ ,  $s > 0$ , that

$$\gamma(st) = \gamma(t)\gamma(s) \quad \text{and} \quad \delta(st) = \gamma(t)\delta(s) + \delta(t).$$

The first equation is the Hamel functional equation, whose only finite, measurable, nonnegative solution is of the following form:

$$\gamma(t) = t^{-\theta}, \quad \theta \in \mathbb{R}.$$

The three cases  $\theta = 0$ ,  $\theta > 0$  and  $\theta < 0$  lead to the Gumbel, Fréchet and Weibull distribution respectively.

For a complete proof see Resnick [2]. □

*Remark 1.2.3.* Usually the Weibull distribution is defined as

$$F(x) = 1 - \exp\left[-\left(\frac{x}{\lambda}\right)^\alpha\right], \quad x \geq 0, \quad \lambda, \alpha > 0$$

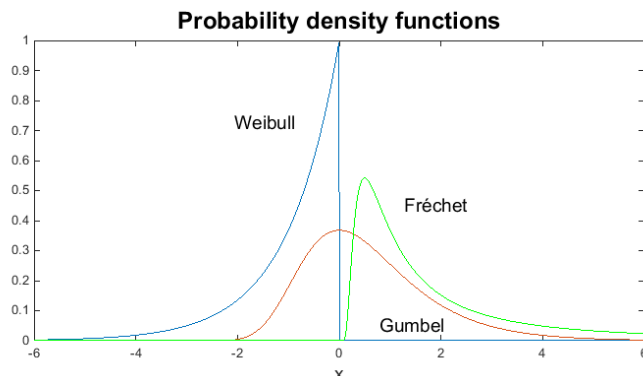
but in the context of extreme value theory is concentrated on  $(-\infty, 0)$ , so

$$\Psi_{\alpha, \lambda} = 1 - F(-x) = \exp\left[-\left(\frac{-x}{\lambda}\right)^\alpha\right], \quad x < 0, \quad \lambda, \alpha > 0.$$

Thus we follow the convention and refers to

$$\Psi_\alpha = \exp(-(-x)^\alpha), \quad x < 0, \quad \alpha > 0$$

as the Weibull distribution (with scale parameter  $\lambda = 1$ ).



**Figure 1.1:** Densities of standard Gumbel (red), Fréchet (green) and Weibull (blue) distributions, with  $\alpha = 1$ .

Even if these families seem to be very different, they are strictly related from a mathematical point of view. Suppose  $X$  be a positive random variable, then the following properties hold:

$$X \sim \Phi_\alpha \Leftrightarrow \ln X^\alpha \sim \Lambda \Leftrightarrow -X^{-1} \sim \Psi_\alpha.$$

Furthermore it is important to note that the limit distribution is unique up to affine transformations: in fact, if the limit appears as

$$\lim_{n \rightarrow \infty} \mathbb{P}(c_n^{-1}(M_n - d_n) \leq x) = H(cx + d),$$

then, by changing norming constants, we obtain

$$\lim_{n \rightarrow \infty} \mathbb{P}(\tilde{c}_n^{-1}(M_n - \tilde{d}_n) \leq x) = H(x)$$

where  $\tilde{c}_n = c_n/c$  and  $\tilde{d}_n = d_n - dc_n/c$ .

This result is due to the convergence to types theorem (see Resnick [2], Proposition 0.2) and is useful to define more general distributions: the *extreme value distributions*, not standardized.

### 1.3 Maximum domain of attraction

In the past section we saw all possible limit distributions of normalised maxima (Theorem 1.2.3); in this section we will show the necessary and sufficient conditions under which a distribution  $F$  gets precisely that limit distribution  $H$ .

**Definition 1.3.1.** We say that the distribution function  $F$  of  $X$  belongs to the *maximum domain of attraction* of the extreme value distribution  $H$  if

$$\exists \quad c_n > 0, d_n \in \mathbb{R} \quad \text{such that} \quad M_n^* := \frac{M_n - d_n}{c_n} \xrightarrow{d} H.$$

In this case we write  $F \in MDA(H)$ .

As direct consequence of Poisson approximation (Theorem 1.2.1) we have

**Proposition 1.3.1.** *The distribution function  $F$  belongs to  $MDA(H)$  with norming constants  $c_n > 0, d_n \in \mathbb{R}$  if and only if*

$$\lim_{n \rightarrow \infty} n\bar{F}(c_n x + d_n) = -\ln H(x), \quad x \in \mathbb{R}.$$

When  $H(x) = 0$  the limit is  $\infty$ .

Another concept which will be useful is the following:

**Definition 1.3.2.** Two distribution functions  $F$  and  $G$  are said to be *tail-equivalent* if

1. they have the same right endpoint (i.e.  $x_F = x_G$ );
2.  $\lim_{x \rightarrow x_F} \bar{F}(x)/\bar{G}(x) = k, \quad 0 < k < \infty$ .

Indeed, it can be proved (see Embrechts et al. [1]) that every maximum domain of attraction is closed to tail-equivalence, that is for two tail-equivalent distribution functions  $F$  and  $G$

$$F \in MDA(H) \quad \Leftrightarrow \quad G \in MDA(H).$$

Before going ahead with the characterization of various types of  $MDA$ , is useful giving the following

**Definition 1.3.3.** A positive, Lebesgue measurable function  $L$  on  $(0, \infty)$  is called

1. *slowly varying* at  $\infty$  ( $L \in \mathcal{R}_0$ ) if

$$\lim_{x \rightarrow \infty} \frac{L(tx)}{L(x)} = 1, \quad t > 0$$



2. *regularly varying* at  $\infty$  of index  $\alpha \in \mathbb{R} \setminus \{0\}$  ( $L \in \mathcal{R}_\alpha$ ) if

$$\lim_{x \rightarrow \infty} \frac{L(tx)}{L(x)} = t^\alpha, \quad t > 0$$

3. *rapidly varying* at  $\infty$  of index  $-\infty$  ( $L \in \mathcal{R}_{-\infty}$ ) if

$$\lim_{x \rightarrow \infty} \frac{L(tx)}{L(x)} = \begin{cases} 0 & \text{if } t > 1 \\ \infty & \text{if } 0 < t < 1 \end{cases}$$

### 1.3.1 Case of Fréchet distribution $\Phi_\alpha$

The tail of Fréchet distribution decreases like a power law: in fact, by Taylor expansion, we obtain

$$\bar{\Phi}_\alpha = 1 - \Phi_\alpha = 1 - \exp(-x^{-\alpha}) = 1 - [1 + (-x^{-\alpha}) + o((-x^{-\alpha})^2)] \sim x^{-\alpha}, \quad x \rightarrow \infty.$$

The following theorem let us know how far away we can move from a power tail and still remain in  $MDA(\Phi_\alpha)$ .

#### Theorem 1.3.2.

*The distribution function  $F$  belongs to  $MDA(\Phi_\alpha)$ ,  $\alpha > 0$ , if and only if*

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

*for some slowly varying function  $L$ .*

*In this case we have  $M_n^* \xrightarrow{d} \Phi_\alpha$  with norming constants  $c_n = F^{\leftarrow}(1 - 1/n)$  and  $d_n = 0$ .*

*Proof.* The proof uses the proposition given above and the Convergence to types theorem. See Embrechts et al. [1] for details.  $\square$

Note that this result implies that every  $F \in MDA(\Phi_\alpha)$  has an infinite right endpoint  $x_F = \infty$ .

This class of distribution functions ( $MDA(\Phi_\alpha)$ ) contains *very heavy-tailed* distributions, that is  $\mathbb{E}[\max(0, X)^\delta] = \infty$  for  $\delta > \alpha$ . For this reason it is appropriate for modelling large insurance claims.

**Example 1.3.1.** *Distributions like Pareto or Cauchy belongs to the MDA of Fréchet distribution, in fact their right tails are of the form*

$$\bar{F}(x) \sim kx^{-\alpha}, \quad x \rightarrow \infty,$$

for some  $k, \alpha > 0$  and this implies  $F \in MDA(\Phi_\alpha)$ .

In particular norming constants are  $d_n = 0$  and  $c_n = (kn)^{1/\alpha}$ .

Also the log-Gamma distribution belongs to  $MDA(\Phi_\alpha)$ , in fact its tail is

$$\bar{F}(x) \sim \frac{\alpha^{\beta-1}}{\Gamma(\beta)} (\ln x)^{\beta-1} x^{-\alpha}, \quad x \rightarrow \infty, \quad \alpha, \beta > 0,$$

where

$$\lim_{x \rightarrow \infty} \frac{(\ln tx)^{\beta-1}}{(\ln x)^{\beta-1}} = \lim_{x \rightarrow \infty} \left( \frac{\ln t}{\ln x} + 1 \right)^{\beta-1} = 1, \quad t > 0,$$

which implies  $(\ln x)^{\beta-1} \in \mathcal{R}_0$  as we wanted.

### 1.3.2 Case of Weibull distribution $\Psi_\alpha$

First is important noting that every  $F \in MDA(\Psi_\alpha)$  has finite right endpoint  $x_F$ . Moreover, since Fréchet and Weibull distributions are related by

$$\Phi_\alpha(x) = \Psi_\alpha(-x^{-1}), \quad x > 0,$$

we can expect that their  $MDA$  are also closely related, as the theorem below shows.

**Theorem 1.3.3.** *The distribution function  $F$  belongs to  $MDA(\Psi_\alpha)$ ,  $\alpha > 0$ , if and only if*

1.  $x_F < \infty$
2.  $\bar{F}(x_F - x^{-1}) = x^{-\alpha} L(x)$  for some slowly varying function  $L$ .

In this case we have  $M_n^* \xrightarrow{d} \Psi_\alpha$  with norming constants  $c_n = x_F - F^{\leftarrow}(1 - 1/n)$  and  $d_n = x_F$ .

*Proof.* The proof is rather formal so we remaind to Embrechts et al. [1] for a sketch. □

Although these distributions have very heavy tails, they're not used for modelling large insurance claims because they are bounded to the right. In fact, even if in practice there is often an upper limit, if we want to use this class of distributions in models we also have to incorporate the parameter  $x_F$  and make models more complicated.

Thus distributions with  $x_F = \infty$  are preferred.

**Example 1.3.2.** Uniform distribution on  $(0, 1)$  belongs to the MDA of Weibull distribution. The right endpoint is  $x_F = 1 < \infty$  and its right tail is of the form

$$\bar{F}(1 - x^{-1}) = x^{-1} \in \mathcal{R}_{-1}$$

which implies  $F \in MDA(\Psi_1)$ .

In particular norming constants are  $c_n = 1 - F^{\leftarrow}(1 - 1/n) = n^{-1}$  and  $d_n = 1$ . More generally can be proved that the Beta distribution belongs to MDA of Weibull distribution.

### 1.3.3 Case of Gumbel distribution $\Lambda$

The tail of Gumbel distribution decreases to 0 at an exponential rate: in fact, by Taylor expansion, we obtain

$$\bar{\Lambda}(x) = 1 - \Lambda(x) = 1 - \exp(-e^{-x}) = 1 - [1 + (-e^{-x}) + o((-e^{-x})^2)] \sim e^{-x}, \quad x \rightarrow \infty.$$

Thus  $MDA(\Lambda)$  contains distribution functions which tails range from *light* (as Normal distribution) to *moderately-heavy* (as log-Normal distribution). Furthermore cases  $x_F = \infty$  and  $x_F < \infty$  are both possible.

**Theorem 1.3.4.** The distribution function  $F$  with right endpoint  $x_F \leq \infty$  belongs to  $MDA(\Lambda)$  if and only if

$$\exists z < x_F \quad \text{such that} \quad \bar{F}(x) = c(x) \exp\left(-\int_z^x \frac{g(t)}{a(t)} dt\right), \quad z < x < x_F$$

where  $c, g$  are measurable functions satisfying  $c(x) \rightarrow c > 0$ ,  $g(x) \rightarrow 1$  as  $x \rightarrow x_F$ , and  $a$  is a positive, absolutely continuous function (respect to Lebesgue measure) with density  $a'(x)$  such that  $\lim_{x \rightarrow x_F} a'(x) = 0$ .

Norming constants can be chosen as  $d_n = F^{\leftarrow}(1 - 1/n)$  and  $c_n = a(d_n)$ .

A possible choice for the function  $a$  is  $a(x) = \int_x^{x_F} \frac{\bar{F}(t)}{\bar{F}(x)} dt$ ,  $x < x_F$ .

*Proof.* The proof is long and technical, see Resnick [2], Corollary 1.7 and Proposition 1.9 for a complete implementation.  $\square$

**Example 1.3.3.** It's rather complicated showing that Normal distribution belongs to MDA of Gumbel distribution because the proof uses the notion of Von Mises function and we didn't introduce it. However, to get an idea we remained to Embrechts et al. [1], Example 3.3.29.

Similarly can be proved that also the Gamma distribution belongs to  $MDA(\Lambda)$ . Using the transformation  $\bar{X} = g(X) = e^{\mu + \sigma x}$ , where  $X \sim N(0, 1)$ ,  $\mu \in \mathbb{R}$

and  $\sigma > 0$ , we have the log-Normal distribution. Since  $X \in MDA(\Lambda)$ , this implies

$$\lim_{n \rightarrow \infty} \mathbb{P}(\tilde{M}_n \leq e^{\mu + \sigma(c_n x + d_n)}) = \Lambda(x), \quad x \in \mathbb{R}$$

where  $c_n$  and  $d_n$  are the norming constant of  $X$ .

This implies

$$\lim_{n \rightarrow \infty} \mathbb{P}(e^{-\mu - \sigma d_n} \tilde{M}_n \leq 1 + \sigma c_n x + o(c_n)) = \Lambda(x), \quad x \in \mathbb{R}.$$

Since  $c_n \rightarrow 0$ , it follows that

$$\frac{1}{\sigma c_n e^{\mu + \sigma d_n}} (\tilde{M}_n - e^{\mu + \sigma d_n}) \xrightarrow{d} \Lambda,$$

so  $\tilde{X} \in MDA(\Lambda)$  with norming constants  $\tilde{c}_n = \sigma c_n e^{\mu + \sigma d_n}$  and  $\tilde{d}_n = e^{\mu + \sigma d_n}$ .

## 1.4 Generalised extreme value distribution

In the past section we saw that the limit distribution of standardised maximum of iid random variables belongs to one of these three families: Fréchet, Weibull, Gumbel.

A one-parameter representation of these distributions can be useful, mainly for statistical applications.

**Definition 1.4.1.** The distribution function defined as

$$H_\xi(x) = \begin{cases} \exp\{-(1 + \xi x)^{-1/\xi}\} & \text{if } \xi \neq 0 \\ \exp\{-e^{-x}\} & \text{if } \xi = 0 \end{cases}$$

where  $1 + \xi x > 0$  and  $H_0$  is obtained as a limit of  $H_\xi$  for  $\xi \rightarrow 0$ , is called *standard generalised extreme value distribution* (SGEV).

The parameter  $\xi$  is called *shape parameter*.

The support of  $H_\xi(x)$  corresponds to

$$x > -1/\xi \quad \text{for } \xi > 0$$

$$x < -1/\xi \quad \text{for } \xi < 0$$

$$x \in \mathbb{R} \quad \text{for } \xi = 0.$$

For different values of  $\xi$  we found the three distribution functions of Fisher-Tippett theorem (1.2.3), in particular

- if  $\xi = 1/\alpha > 0$  we have the Fréchet distribution  $\Phi_\alpha$ ;
- if  $\xi = -1/\alpha < 0$  we have the Weibull distribution  $\Psi_\alpha$ ;
- if  $\xi = 0$  we have the Gumbel distribution  $\Lambda$ .

Consequently, we can also give a unique characterisation of the maximum domain of attraction which includes characteristics of each distribution family.

**Theorem 1.4.1.** *The distribution function  $F$  belongs to  $MDA(H_\xi)$ ,  $\xi \in \mathbb{R}$ , if and only if there exists a positive measurable function  $a$  such that for  $1 + \xi x > 0$*

$$\lim_{u \rightarrow x_F} \frac{\bar{F}(u + a(u)x)}{\bar{F}(u)} = \begin{cases} (1 + \xi x)^{-1/\xi} & \text{if } \xi \neq 0 \\ e^{-x} & \text{if } \xi = 0. \end{cases}$$

*Proof.* We only sketch the case  $\xi > 0$  reminding to Embrechts et al. [1], Theorem 3.4.5 for more details.

( $\Rightarrow$ ) For  $\xi > 0$  we have  $H_\xi(x) = \Phi_\alpha((x + \alpha)/\alpha)$ ,  $\alpha = 1/\xi$ , and for Theorem (1.3.2)  $F \in MDA(H_\xi)$  is equivalent to  $\bar{F} \in \mathcal{R}_{-\alpha}$ . Using the representation theorem for regularly varying functions we obtain

$$\lim_{u \rightarrow \infty} \frac{\bar{F}(u + a(u)x)}{\bar{F}(u)} = \left(1 + \frac{x}{\alpha}\right)^{-\alpha},$$

which is the relation above.

( $\Leftarrow$ ) If the relation holds, chosen  $d_n = (1/\bar{F})^{\leftarrow}(n)$ , then  $1/\bar{F}(d_n) \sim n$ . Substituting  $u = d_n$  in the relation above, we obtain

$$\left(1 + \frac{x}{\alpha}\right)^{-\alpha} = \lim_{n \rightarrow \infty} \frac{\bar{F}(d_n + a(d_n)x)}{\bar{F}(d_n)} = \lim_{n \rightarrow \infty} n \bar{F}(d_n + a(d_n)x),$$

and for Proposition (1.3.1)  $F \in MDA(H_\xi)$ ,  $\xi = 1/\alpha$ . □

*Remark 1.4.1.* We can also introduce the location-scale family  $H_{\xi;\mu,\sigma}$ , with  $\mu \in \mathbb{R}$  and  $\sigma > 0$ , by replacing the argument  $x$  above with  $(x - \mu)/\sigma$ . We denote it with GEV and it's support has to be adjusted accordingly.

In particular we obtain

$$H_{\xi;\mu,\sigma}(x) = \begin{cases} \exp \left\{ - \left[ 1 + \xi \left( \frac{x-\mu}{\sigma} \right) \right]^{-1/\xi} \right\} & \text{if } \xi \neq 0 \\ \exp \left\{ - e^{-\left( \frac{x-\mu}{\sigma} \right)} \right\} & \text{if } \xi = 0 \end{cases}$$

which support is  $x > \mu - \sigma/\xi$  if  $\xi > 0$ ,  $x < \mu - \sigma/\xi$  if  $\xi < 0$  and  $x \in \mathbb{R}$  if  $\xi = 0$ .

Mean and variance for SGEV distribution are obtained using moments (see Hosking and Wallis [3], Kotz and Nadarajah [4], Reiss and Thomas [5] p. 17-18):

$$m_j = \begin{cases} (-1)^j \Gamma(1 - j\xi) & \text{if } \xi < 0 \\ \Gamma(1 - j\xi) & \text{if } \xi \in (0, 1/j) \\ \infty & \text{if } \xi \geq 1/j. \end{cases}$$

In particular if  $X$  has  $H_\xi$  distribution function we have

$$\mathbb{E}[X] = m_1 = \begin{cases} -\Gamma(1 - \xi) & \text{if } \xi < 0 \\ \Gamma(1 - \xi) & \text{if } \xi \in (0, 1) \\ \infty & \text{if } \xi \geq 1 \end{cases}$$

$$Var(X) = m_2 - m_1^2 = \Gamma(1 - 2\xi) - [\pm\Gamma(1 - \xi)]^2 \quad \text{if } \xi < 1/2, \quad \xi \neq 0.$$

*Remark 1.4.2.* For a variable  $X$  with  $H_{\xi,\mu,\sigma}$  distribution function we have

$$\mathbb{E}[X] = \begin{cases} \mu - \frac{\sigma}{\xi}(1 - \Gamma(1 - \xi)) & \text{if } \xi < 1 \text{ and } \xi \neq 0 \\ \infty & \text{if } \xi \geq 1 \end{cases}$$

$$Var(X) = \begin{cases} \frac{\sigma^2}{\xi^2}(\Gamma(1 - 2\xi) - [\Gamma(1 - \xi)]^2) & \text{if } \xi < 1/2 \text{ and } \xi \neq 0 \\ \infty & \text{if } \xi \geq 1/2. \end{cases}$$

If  $\xi = 0$ , using the moment generating function

$$\mathbb{E}[e^{tx}] = e^{t\mu} \Gamma(1 - \sigma t), \quad \sigma|t| < 1,$$

we have

$$\mathbb{E}[X] = \mu + \gamma\sigma$$

where  $\gamma = 0.57722$  is the Euler's constant, and

$$Var(X) = \frac{\pi^2}{6}\sigma^2.$$

Quantiles of  $H_\xi$  are given by

$$x_p = \begin{cases} [(-\ln p)^{-\xi} - 1]/\xi & \text{if } \xi \neq 0 \\ -\ln(-\ln p) & \text{if } \xi = 0 \end{cases}, 0 < p < 1$$

and they have a special meaning:  $x_p = H_\xi^{\leftarrow}(p)$  is the **return level** associated with the **return period**  $1/(1-p)$ , that is the level  $x_p$  is expected to be exceeded on average once every  $1/(1-p)$  years. The quantity  $1-p$  is called **probability of exceedance** and represents the probability that  $M_n^*$  falls beyond the threshold  $x_p$ .

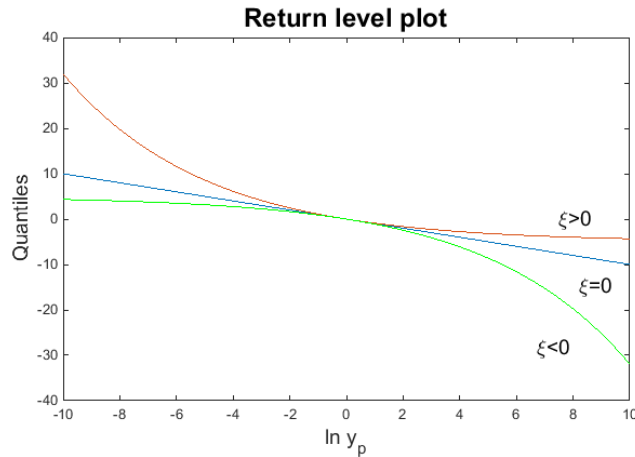
In particular, substituting  $-\ln p$  with  $y_p$ , we obtain

$$x_p = \begin{cases} (y_p^{-\xi} - 1)/\xi & \text{if } \xi \neq 0 \\ -\ln y_p & \text{if } \xi = 0 \end{cases}$$

and plotting quantiles against  $\ln y_p =: w_p$  ( $w_p \rightarrow -\infty$  as  $1-p \rightarrow 0$ ),

$$x_p = \begin{cases} (e^{-w_p \xi} - 1)/\xi & \text{if } \xi \neq 0 \\ -w_p & \text{if } \xi = 0 \end{cases}$$

we can see that the plot is linear when  $\xi = 0$ , convex with no finite bound for  $1-p \rightarrow 0$  when  $\xi > 0$  and concave with asymptotic limit  $x_p = -1/\xi$  for  $1-p \rightarrow 0$  when  $\xi < 0$ . This kind of graph is called **return level plot**.



**Figure 1.2:** Return level plots of the SGEV distribution with  $\xi = 0.2$ ,  $\xi = -0.2$  and  $\xi = 0$ .

## 1.5 Inference for GEV distribution

In this section we refer to  $H_{\xi;\mu,\sigma}$ , with  $\xi, \mu \in \mathbb{R}$ ,  $\sigma > 0$ , and our purpose is showing methods used for estimating GEV distribution parameters.

Since GEV distribution is the limit distribution of the normalised maxima of a sequence of iid random variables, for estimating its parameters is necessary having a sequence of independent variables all with the same GEV distribution function of the normalised maxima.

For this purpose we consider a sequence of independent random variables  $X_1, \dots, X_{nm}$ , for large value of  $n$ , with the same distribution function  $F$  and we block them in  $m$  groups of equal length  $n$ ; group's maxima are called **block maxima** and are denoted by  $Z_i$ ,  $i = 1, \dots, m$ .

Notice that the choice of  $n$  is critical: in fact if  $n$  is too small there would be bias in estimation and extrapolation due to poor approximation of the limit distribution given in theorem 1.2.3; on the other hand if  $n$  is too big there would be few groups, leading to large estimation variance.

Practical considerations often lead to blocks of length one year: in this way is plausible assuming that block maxima have a common distribution.

Furthermore, is useful also that  $Z_1, Z_2, \dots, Z_m$  are independent so that we can easily write the log-likelihood.

*Remark 1.5.1.* If the  $X_j$ ,  $j = 1, \dots, nm$ , are independent the  $Z_i$ ,  $i = 1, \dots, m$  are also independent. However, independence of the  $Z_i$ 's is likely to be a reasonable approximation even if the  $X_j$ 's constitute a dependent series (see Coles [6]).

The use of these block maxima for the estimation of  $M_{nm}^*$  distribution parameters is justified as follows.

The idea is to consider  $M_{nm}^*$  as maximum of  $nm$  variables and, at the same time, as maximum of  $m$  maxima, each of which is the maximum of  $n$  variables. More precisely we know from theorem 1.2.3 that, for  $n$  large,

$$\mathbb{P}(M_n^* \leq x) = \mathbb{P}\left(\frac{M_n - d_n}{c_n} \leq x\right) \approx H_{\xi;\mu,\sigma}(x),$$

so for any  $m \in \mathbb{N}$ , since  $nm$  is also large, we have

$$\mathbb{P}(M_{nm}^* \leq x) = \mathbb{P}\left(\frac{M_{nm} - d_{nm}}{c_{nm}} \leq x\right) \approx H_{\xi;\mu,\sigma}(x).$$



On the other hand, being  $M_{nm}^*$  the maximum of  $m$  maxima each of which is the maximum of  $n$  variables, we obtain

$$\mathbb{P}(M_{nm}^* \leq x) = \left[ \mathbb{P}\left(\frac{M_n - d_n}{c_n} \leq x\right) \right]^m \approx H_{\xi;\mu,\sigma}^m(x).$$

Therefore  $H_{\xi;\mu,\sigma}$  and  $H_{\xi;\mu,\sigma}^m$  are identical apart from norming constants used on  $M_{nm}$ .

Then let us consider  $Z_i$ ,  $i = 1, \dots, m$  as independent random variables with GEV distribution

$$H_{\xi;\mu,\sigma}(z_i) = \begin{cases} \exp\left\{-\left[1 + \xi\left(\frac{z_i - \mu}{\sigma}\right)\right]^{-1/\xi}\right\} & \text{if } \xi \neq 0 \\ \exp\left\{-e^{-\left(\frac{z_i - \mu}{\sigma}\right)}\right\} & \text{if } \xi = 0 \end{cases} \quad (1.2)$$

where  $1 + \xi\left(\frac{z_i - \mu}{\sigma}\right) > 0$ , whose density functions are given by

$$h_{\xi;\mu,\sigma}(z_i) = \begin{cases} \left[1 + \xi\left(\frac{z_i - \mu}{\sigma}\right)\right]^{-1-1/\xi} H_{\xi;\mu,\sigma}(z_i)/\sigma & \text{if } \xi \neq 0 \\ e^{-\left(\frac{z_i - \mu}{\sigma}\right)} H_{\xi;\mu,\sigma}(z_i)/\sigma & \text{if } \xi = 0. \end{cases}$$

The log-likelihood is

$$\begin{aligned} l(\mu, \sigma, \xi) &= \ln \left[ \prod_{i=1}^m h_{\xi;\mu,\sigma}(z_i) \right] = \\ &= \begin{cases} -m \ln \sigma - \left(1 + \frac{1}{\xi}\right) \sum_{i=1}^m \ln \left[1 + \xi\left(\frac{z_i - \mu}{\sigma}\right)\right] - \sum_{i=1}^m \left[1 + \xi\left(\frac{z_i - \mu}{\sigma}\right)\right]^{-1/\xi} & \text{if } \xi \neq 0 \\ -m \ln \sigma - \sum_{i=1}^m \left(\frac{z_i - \mu}{\sigma}\right) - \sum_{i=1}^m e^{-\left(\frac{z_i - \mu}{\sigma}\right)} & \text{if } \xi = 0 \end{cases} \end{aligned} \quad (1.3)$$

and their maximization with respect the parameter vector  $(\mu, \sigma, \xi)$ , obtained equating partial derivatives of equations (1.3) to zero, gives the maximum likelihood estimates  $\hat{\mu}$ ,  $\hat{\sigma}$  and  $\hat{\xi}$ ; there is no analytical solution, but it can be achieved using standard numerical methods.

The estimates vector  $(\hat{\mu}, \hat{\sigma}, \hat{\xi})$  has approximately a multivariate normal distribution with mean  $(\mu, \sigma, \xi)$ , so it's an unbiased estimator for parameters of GEV distribution.

Quantiles estimates can then be obtained substituting  $(\hat{\mu}, \hat{\sigma}, \hat{\xi})$  in the quantile expression for GEV distribution with location and scale parameters, leading to

$$\hat{z}_p = \begin{cases} \hat{\mu} + [(-\ln p)^{-\hat{\xi}} - 1] \hat{\sigma} / \hat{\xi} & \text{if } \hat{\xi} \neq 0 \\ \hat{\mu} - \hat{\sigma} \ln(-\ln p) & \text{if } \hat{\xi} = 0. \end{cases} \quad (1.4)$$

The variance-covariance matrix of quantiles estimates is calculated by the delta method, which we briefly illustrate.

Let us consider a real valued continuous and differentiable function  $h$  of a consistent estimator  $\hat{\boldsymbol{\theta}}$  of  $\boldsymbol{\theta}$ , defined in a neighbourhood of  $\boldsymbol{\theta}$ , such that  $\nabla h(\boldsymbol{\theta}) \neq \mathbf{0}$  and write its first-order Taylor expansion

$$h(\hat{\boldsymbol{\theta}}) \sim h(\boldsymbol{\theta}) + \nabla h(\boldsymbol{\theta})^T (\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}).$$

Then

$$\begin{aligned} \text{Var}(h(\hat{\boldsymbol{\theta}})) &= \text{Var}(h(\boldsymbol{\theta}) + \nabla h(\boldsymbol{\theta})^T (\hat{\boldsymbol{\theta}} - \boldsymbol{\theta})) \\ &= \text{Var}(h(\boldsymbol{\theta}) + \nabla h(\boldsymbol{\theta})^T \hat{\boldsymbol{\theta}} - \nabla h(\boldsymbol{\theta})^T \boldsymbol{\theta}) \\ &= \text{Var}(\nabla h(\boldsymbol{\theta})^T \hat{\boldsymbol{\theta}}) \\ &= \nabla h(\boldsymbol{\theta})^T \text{Var}(\hat{\boldsymbol{\theta}}) \nabla h(\boldsymbol{\theta}). \end{aligned}$$

For more details see Davison [7].

Since  $\hat{z}_p$  is a function of  $(\hat{\mu}, \hat{\sigma}, \hat{\xi})$ , we can obtain

$$\text{Var}(\hat{z}_p) \sim \nabla z_p^T V \nabla z_p$$

where  $V$  is the variance-covariance matrix of  $(\hat{\mu}, \hat{\sigma}, \hat{\xi})$  and

$$\nabla z_p = \begin{bmatrix} \frac{\partial z_p}{\partial \mu} \\ \frac{\partial z_p}{\partial \sigma} \\ \frac{\partial z_p}{\partial \xi} \end{bmatrix} = \begin{cases} \begin{bmatrix} 1 \\ [(-\ln p)^{-\xi} - 1]/\xi \\ \{-(-\ln p)^{-\xi}[\xi \ln(-\ln p) + 1] + 1\}\sigma/\xi^2 \end{bmatrix} & \text{if } \xi \neq 0 \\ \begin{bmatrix} 1 \\ -\ln(-\ln p) \\ 0 \end{bmatrix} & \text{if } \xi = 0. \end{cases}$$

In the case in which  $\hat{\xi} < 0$  it's also possible to give the maximum likelihood estimate of the right endpoint  $z_F$ , because for the Weibull distribution it is finite. Since  $z_F$  corresponds to  $z_p$  where  $p = 1$ , we have

$$\hat{z}_F = \hat{\mu} - \hat{\sigma}/\hat{\xi},$$

where its variance-covariance matrix is evaluated for  $p = 1$ .

Another method for making inference on a particular component of the parameter vector  $(\mu, \sigma, \xi)$  is the profile likelihood. It consists in fixing the parameter of interest and maximizing the log-likelihood (1.3) with respect to the remaining parameters. This is repeated for a range of values of the parameter of interest. The corresponding maximized values of the log-likelihood constitute the profile likelihood, used to obtain approximate confidence intervals (see Coles [6]).

Other methods can then be used for model checking:

1. the *probability plot* is a comparison of the empirical and fitted distribution functions.

With ordered block maximum data  $z_{(1)} \leq \dots \leq z_{(m)}$ , the empirical distribution function is given by  $\tilde{H}(z_{(i)}) = i/(m+1)$  and the fitted distribution function  $\hat{H}(z_{(i)})$  is obtained by substitution of parameter estimates into expression (1.2) with  $z_i = z_{(i)}$ ,  $i = 1, \dots, m$ .

If the GEV model works well  $\tilde{H}(z_{(i)}) \sim \hat{H}(z_{(i)})$  and the probability plot, consisting of the locus of points

$$\left\{ \left( \tilde{H}(z_{(i)}), \hat{H}(z_{(i)}) \right), i = 1, \dots, m \right\}$$

should lie close to the unit diagonal. Its weakness for extreme value models is that both distributions tend toward approach to 1 as  $z_{(i)}$  increases, so probability plot provides the least information in the region of more interest;

2. the *quantile plot* solves the weakness of probability plot because it consists of the points

$$\left\{ \left( \hat{H}^{-1}(i/(m+1)), z_{(i)} \right), i = 1, \dots, m \right\}.$$

Departures from linearity indicate model failure;

3. the *return level plot* consists of the locus of points

$$\{(\ln y_p, \hat{z}_p) : 0 < p < 1\}$$

and is particularly convenient in extreme value models because the tail of the distribution is compressed and return level estimates  $\hat{z}_p$  for long return periods can be displayed. Furthermore, the linearity of the plot in the case  $\xi = 0$  is useful in judging the effect of the estimated shape parameter. Empirical estimates of the return level function found above can also be added in order to use return level plot as model diagnostic;

4. another equivalent diagnostic, not based on comparison of empirical and fitted distribution functions, consists in comparing the probability density function of the fitted model with the histogram of the data.

## 1.6 Threshold models

Another point of view used for modeling extremal events, especially when other data on extremes are available, is considering events which exceed a certain threshold. Theorem 1.4.1 gives an interesting interpretation.

In fact, reformulating the relation found there, we have

$$\lim_{u \rightarrow x_F} \mathbb{P} \left( \frac{X - u}{a(u)} > x \mid X > u \right) = \begin{cases} (1 + \xi x)^{-1/\xi} & \text{if } \xi \neq 0 \\ e^{-x} & \text{if } \xi = 0 \end{cases}$$

where  $X$  is a random variable with distribution function  $F \in MDA(H_\xi)$ .

This relation represents the limit distribution for the *scaled excess* over the (high) threshold  $u$ , with scaling factor  $a(u)$ .

**Definition 1.6.1.** Let  $X$  be a random variable with distribution function  $F$  and upper endpoint  $x_F$ . For fixed  $u < x_F$ , the function

$$F_u(x) = \mathbb{P}(X - u \leq x \mid X > u), \quad x \geq 0,$$

is called the *excess distribution function* of the random variable  $X$  over the threshold  $u$ . The function

$$e(u) = \mathbb{E}[X - u \mid X > u]$$

is called the *mean excess function* of  $X$ .

*Remark 1.6.1.* In insurance context the function  $F_u$  is usually called *excess of loss distribution function*.

As we can expect, the complement to 1 of the reformulation at the beginning of this section is a limit distribution function, in particular we have the following

**Definition 1.6.2.** Let  $X$  be a random variable with distribution function  $F \in MDA(H_\xi)$ . We call *standard generalised Pareto distribution* (SGPD) the distribution function

$$G_\xi(x) = \begin{cases} 1 - (1 + \xi x)^{-1/\xi} & \text{if } \xi \neq 0 \\ 1 - e^{-x} & \text{if } \xi = 0, \end{cases}$$

where  $G_0$  can be interpreted as the limit of  $G_\xi$  for  $\xi \rightarrow 0$  and the support is

$$\begin{aligned} x &\geq 0 \quad \text{if } \xi \geq 0 \\ 0 \leq x &\leq -1/\xi \quad \text{if } \xi < 0. \end{aligned}$$

Therefore the SGPD represents the limit distribution function of the scaled excess over the threshold  $u$ :

$$\left( \frac{X - u}{a(u)} \middle| X > u \right) \xrightarrow{d} G_\xi.$$

*Remark 1.6.2.* We can also introduce the location-scale family  $G_{\xi, \tilde{\mu}, \tilde{\sigma}}$ , with  $\tilde{\mu} \in \mathbb{R}$  and  $\tilde{\sigma} > 0$ , by replacing the argument  $x$  above with  $(x - \tilde{\mu})/\tilde{\sigma}$ . We denote it with GPD and it's support has to be adjusted accordingly.

In particular we obtain

$$G_{\xi, \tilde{\mu}, \tilde{\sigma}}(x) = \begin{cases} 1 - (1 + \xi \frac{x - \tilde{\mu}}{\tilde{\sigma}})^{-1/\xi} & \text{if } \xi \neq 0 \\ 1 - e^{-\frac{x - \tilde{\mu}}{\tilde{\sigma}}} & \text{if } \xi = 0, \end{cases}$$

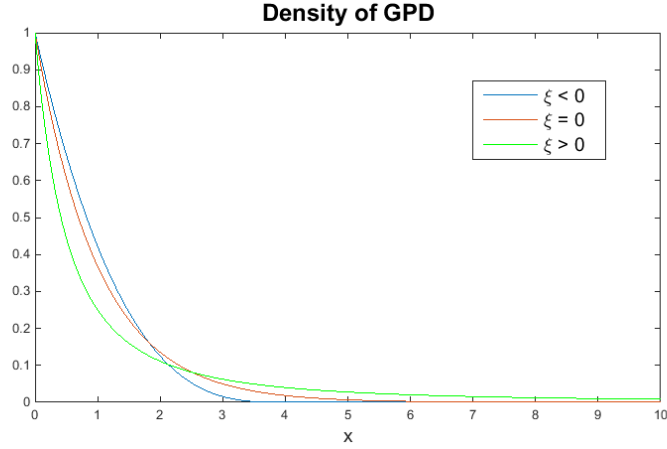
which support is  $x \geq \tilde{\mu}$  if  $\xi \geq 0$ ,  $\tilde{\mu} \leq x \leq \tilde{\mu} - \tilde{\sigma}/\xi$  if  $\xi < 0$ .

Mean and variance of a random variable  $X$  with  $G_{\xi, \tilde{\mu}, \tilde{\sigma}}$  distribution function are given by (see Suzuki-Parker [8])

$$\begin{aligned} \mathbb{E}[X] &= \tilde{\mu} + \frac{\tilde{\sigma}}{1 - \xi} \quad \text{if } \xi < 1 \\ \text{Var}(X) &= \frac{\tilde{\sigma}^2}{(1 - \xi)^2(1 - 2\xi)} \quad \text{if } \xi < \frac{1}{2}. \end{aligned}$$

Duality between the GEV distribution and the GPD implies that the shape parameter  $\xi$  is the same for both distributions, in particular it establishes the qualitative behaviour of generalised Pareto distribution. As we can expect, the generalized Pareto distribution has three basic forms, each corresponding to a limiting distribution of excess data from a different class of underlying distributions:

- distributions whose tails decrease exponentially, such as the Normal, lead to a generalized Pareto shape parameter  $\xi = 0$ ;
- distributions whose tails decrease as a polynomial, such as the log-Gamma, lead to a positive shape parameter  $\xi > 0$ ;



**Figure 1.3:** Densities of generalised Pareto distribution for  $\tilde{\mu} = 0$ ,  $\tilde{\sigma} = 1$  and different values of  $\xi$ :  $\xi = -0.25$  (blue),  $\xi = 0$  (red),  $\xi = 1$  (green).

- distributions whose tails are finite, such as the Beta, lead to a negative shape parameter  $\xi < 0$ .

The generalised Pareto distribution has some interesting properties. The most important is that

$$\forall \xi \in \mathbb{R}, \quad F \in MDA(H_\xi) \Leftrightarrow \lim_{u \rightarrow x_F} \sup_{0 < x < x_F - u} \|F_u(x) - G_{\xi, \tilde{\sigma}(u)}(x)\| = 0$$

for some positive function  $\tilde{\sigma}(u)$ .

As a consequence of this result we will assume  $\tilde{\mu} = 0$ , as in many statistical applications.

Let us show now these properties:

- the class of GPDs is closed with respect to changes of the threshold, that is

$$\frac{\bar{G}_{\xi; \tilde{\sigma}}(x_1 + x_2)}{\bar{G}_{\xi; \tilde{\sigma}}(x_1)} = \bar{G}_{\xi; \tilde{\sigma} + \xi x_1}(x_2) = \left(1 + \xi \frac{x_2}{\tilde{\sigma} + \xi x_1}\right)^{-1/\xi}$$

$$\text{for } x_1, x_2 \in \begin{cases} [0, \infty) & \text{if } \xi \geq 0 \\ [0, -\tilde{\sigma}/\xi] & \text{if } \xi < 0 \end{cases};$$

- if  $X$  has GPD with  $\xi < 1$ , then for  $u < x_F$

$$e(u) = \mathbb{E}[X - u | X > u] = \frac{\tilde{\sigma} + \xi u}{1 - \xi}, \quad \tilde{\sigma} + \xi u > 0.$$

This implies that, given an iid sample  $x_1, \dots, x_n$ , the range of thresholds we can choose is the set of  $u$  for which the empirical mean excess function  $e_n(u)$  is roughly linear. This graph is called **mean residual life plot**.

In the case  $\xi \geq 1$  the mean is infinite;

- if in a model the number of exceedances  $N$  is exactly Poisson and the excess distribution function is an exact GPD, then the maximum of these excesses has an exact GEV distribution.

In other words, if  $N \sim \mathcal{P}(\lambda)$  is independent of the independent and identically distributed sequence  $X_1, \dots, X_N$  with a GPD with parameters  $\xi$  and  $\tilde{\sigma}$ , then the maximum  $M_N = \max\{X_1, \dots, X_N\}$  is such that

$$\mathbb{P}(M_N \leq x) = \exp \left[ -\lambda \left( 1 + \xi \frac{x}{\tilde{\sigma}} \right)^{-1/\xi} \right] = H_{\xi; \mu, \sigma}(x)$$

where  $\mu = \tilde{\sigma}(\lambda^\xi - 1)/\xi$  and  $\sigma = \tilde{\sigma}\lambda^\xi$ .

*Proof.* We know that  $\mathbb{P}(M_n \leq x) = G_{\xi; \tilde{\sigma}}^n(x)$  if  $n$  is fixed and, from remark (1.2.2), that the number of excesses is roughly Poisson. Then we have

$$\begin{aligned} \mathbb{P}(M_N \leq x) &= \sum_{n=0}^{\infty} e^{-\lambda} \frac{\lambda^n}{n!} G_{\xi; \tilde{\sigma}}^n(x) \\ &= e^{-\lambda} \sum_{n=0}^{\infty} \frac{[\lambda G_{\xi; \tilde{\sigma}}(x)]^n}{n!} \\ &= e^{-\lambda} e^{\lambda G_{\xi; \tilde{\sigma}}(x)} \\ &= e^{-\lambda} e^{\lambda(1 - \bar{G}_{\xi; \tilde{\sigma}}(x))} \\ &= e^{-\lambda \bar{G}_{\xi; \tilde{\sigma}}(x)} \\ &= \exp \left[ -\lambda \left( 1 + \xi \frac{x}{\tilde{\sigma}} \right)^{-1/\xi} \right] \\ &= \exp \left[ -\left( \lambda^{-\xi} \left( 1 + \xi \frac{x}{\tilde{\sigma}} \right) \right)^{-1/\xi} \right] \end{aligned}$$

$$\begin{aligned}
&= \exp \left[ - \left( 1 - 1 + \frac{1}{\lambda^\xi} + \xi \frac{x}{\tilde{\sigma} \lambda^\xi} \right) \right]^{-1/\xi} \\
&= \exp \left[ - \left( 1 + \xi \frac{x - \tilde{\sigma}(\lambda^\xi - 1)/\xi}{\tilde{\sigma} \lambda^\xi} \right) \right]^{-1/\xi} \\
&= H_{\xi; \tilde{\sigma}(\lambda^\xi - 1)/\xi, \tilde{\sigma} \lambda^\xi}(x).
\end{aligned}$$

The case  $\xi = 0$  reduces to

$$\mathbb{P}(M_N \leq x) = \exp(-e^{-(x - \tilde{\sigma} \ln \lambda)/\tilde{\sigma}});$$

□

## 1.7 Inference for threshold models

Let us consider an iid sample  $x_1, \dots, x_n$ . Extreme events are those exceeding the threshold  $u$ , that is  $\{x_i | x_i > u, \quad i = 1, \dots, k\}$ . Ordering these exceedances  $x_{(1)} < \dots < x_{(k)}$ , we can define threshold excess as

$$y_i = x_{(i)} - u, \quad i = 1, \dots, k$$

which have the generalised Pareto distribution, not standardised. Assuming  $\tilde{\mu} = 0$ , the  $y_i$ 's densities for  $i = 1, \dots, k$  are of the form

$$g_{\xi; \tilde{\sigma}}(y_i) = \begin{cases} \frac{1}{\tilde{\sigma}} \left( 1 + \xi \frac{y_i}{\tilde{\sigma}} \right)^{-1-1/\xi} & \text{if } \xi \neq 0 \\ \frac{1}{\tilde{\sigma}} e^{-\frac{y_i}{\tilde{\sigma}}} & \text{if } \xi = 0. \end{cases}$$

*Remark 1.7.1.* As for the choice of blocks dimension, also in the threshold choice attention is required. In fact, if the threshold is

- too low the asymptotic basis of the model would be violated leading to bias;
- too high there would be few excesses leading to high variance.

As for the GEV distribution, inference on generalised Pareto distribution consists on maximum likelihood estimation of the parameters vector  $(\tilde{\sigma}, \xi)$ .



So parameters estimates are obtained equating to zero the partial derivatives of the log-likelihood, which is given by the following formula

$$\begin{aligned}
 l(\tilde{\sigma}, \xi) &= \ln \left[ \prod_{i=1}^k g_{\xi; \tilde{\sigma}}(y_i) \right] = \\
 &= \begin{cases} -k \ln \tilde{\sigma} - \left(1 + \frac{1}{\xi}\right) \sum_{i=1}^k \ln \left[1 + \xi \frac{y_i}{\tilde{\sigma}}\right] & \text{if } \xi \neq 0 \\ -k \ln \tilde{\sigma} - \frac{1}{\tilde{\sigma}} \sum_{i=1}^k y_i & \text{if } \xi = 0. \end{cases} \quad (1.5)
 \end{aligned}$$

Also in this case no analytical solution exists, so numerical methods are required.

For what concerns model checking, the fitted generalised Pareto model can be checked using probability plots, quantile plots, return level plots and density plots as explained in Section (1.5).

## CHAPTER 2

# STATISTICAL ANALYSIS OF ITALIAN EARTHQUAKES DATA

### 2.1 Introduction

The theory of extreme events developed in the past chapter can now be applied to Italian earthquakes data.

In this chapter we will analyze data from the last Parametric Catalogue of Italian Earthquakes (CPTI15), in which is recorded information on events from year 1000 to 2014. In addition to general parameters, such as origin time and zone, also macroseismic and instrumental parameters are available:

- macroseismic parameters concern damage effects due to the passage of seismic waves on urban centres, single buildings and on people. These parameters are the result of an elaboration of information collected in the field by teams of experts (i.e. seism classification by Mercalli scale);
- instrumental parameters are those obtainable from instruments, that is epicentre location and magnitude.

Is important underlying the difference between Mercalli and Richter scale. The first was introduced by Giuseppe Mercalli (1850-1914) with the purpose of measuring the earthquake intensity from destructive effects on buildings and people: it's a discrete scale and consists of 12 levels. The second one, introduced by Charles Francis Richter (1900-1985), measures the energy released by an earthquake at the point of fracture of the Earth's crust, named

focus. This is a continuous scale based on the "magnitude", a dimensionless quantity.

<b>Richter scale</b> ( <i>magnitude</i> )	<b>Mercalli scale</b> ( <i>degree</i> )	<b>Perception level</b>
0	I	instrumental
1	I	instrumental
2	I-II	instrumental/feeble
3	III-IV	slight/moderate
4	V	rather strong
5	VI-VII	strong/very strong
6	VIII-IX	rouinous/disastrous
7	X-XI	very disastrous/catastrophic
8	XII	apocalyptic
9	XII	apocalyptic

Because of its scarce relation with physical characteristics of earthquake cause, this scale was recently substituted by the Moment Magnitude scale, introduced by H. Kanamori and T.C. Hanks in 1979, defined by

$$M_w = \frac{2}{3}(\log_{10}M_0 - 6.03)$$

where  $M_0$  is the seismic moment at the focus measured in  $N \cdot m$  and constants are chosen such that having values similar to Richter scale.

The energy released by an earthquake, strictly related with it's destroying power, is proportional (for less then a constant) to oscillation width raised to the 3/2 power. So, in term of released energy, a magnitudo difference of 1 is equivalent to a factor of  $10^{1 \cdot 3/2} = 31.6$ , a magnitudo difference of 2 is equivalent to a factor of  $10^{2 \cdot 3/2} = 1000$  and in general, a magnitudo difference of  $m_2 - m_1$  is equivalent to a factor of

$$f_{\Delta E} \simeq 10^{\frac{3}{2}(m_2 - m_1)}$$

(see Kanamori [13]).

Other information in the Parametric Catalogue is about the *completeness*. Since the period under study is too large, our knowledge of historical earthquakes derives from archive data which often are fragmentary and not completely reliable. This feature intensifies moving across the Italian territory: in fact floods, fires and other located events may have destroyed

archives, causing a loss of information. For this purpose each record presents an item related to the year from which we can consider information satisfactory and reliable.

Finally, we can also distinguish main shocks from replicates.

## 2.2 Preparation of the dataset

First of all, in order to develop our analysis, we need cleaning the dataset<sup>1</sup>. The entire Catalogue consists of 4390 records belonging both to Italian territory (dry land and sea) and neighboring foreign Countries. The reason of this can be immediately found keeping in mind that effects of an earthquake whose epicentre is located not far from the border with Italy can be felt even in Italian territory. Note that foreign data in the dataset were kept directly from foreign catalogues for more accuracy.

Observations not in the completeness period must be dropped, and also records representing replicates: in fact the theory in the first chapter is valid for independent events only.

After this procedure our dataset consists of 1728 records.

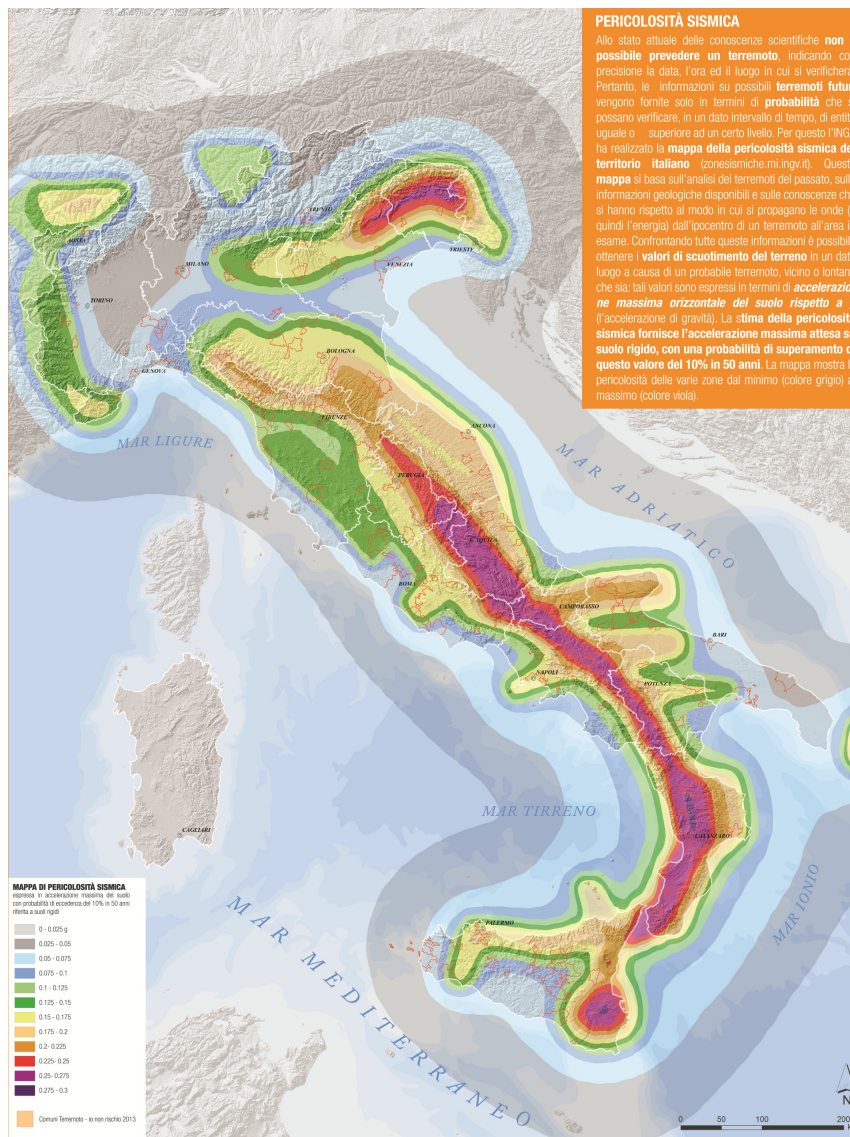
The variables used in this analysis are the following:

1. **Year** is the year of earthquake origin;
2. **MwDef** is the default moment magnitude, that is the weighted average of macroseismic and instrumental moment magnitude with weights the square inverse of their respectively errors;
3. **region** is the completeness zone in which is located the epicentre. We distinguish 6 areas:
  - *Sea/Foreign*
  - *Alps*
  - *Po valley*
  - *Centre*
  - *South*
  - *Islands*.

---

<sup>1</sup>For the preparation of the dataset and all elaborations realised we had used Stata software.

One can ask why we consider the variable **region** such defined. The reason derives from a geological cause. Italian territory in fact is distinguished from a seismic point of view by the presence of faults, which obviously are not homogeneously distributed. An idea can be given by the following map.

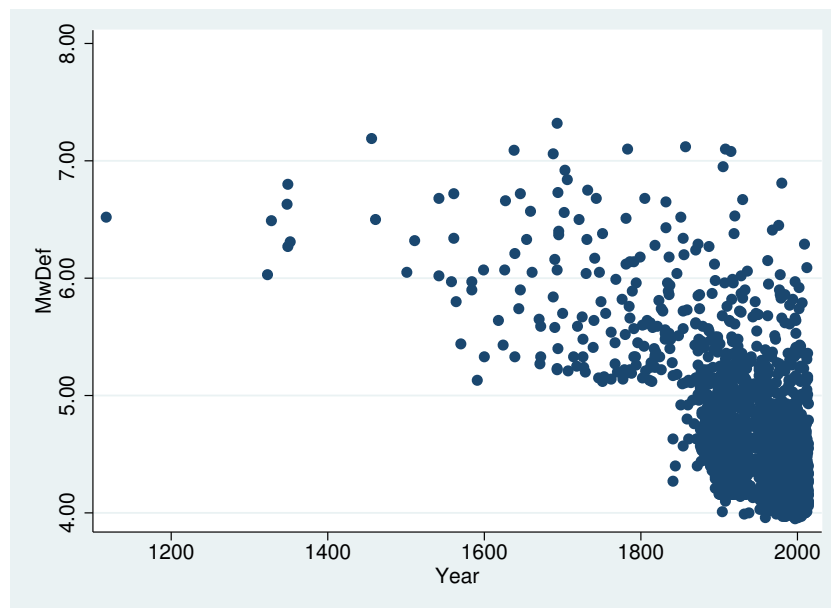


**Figure 2.1:** Seismic hazard map (2013). Violet areas are the most seismic, grey's the less ones. For the source see ING [18].

## 2.3 GEV distribution parameters estimation

### 2.3.1 Global analysis

Consider the whole set of Italian records, consisting of 1728 independent observations. We can call them  $x_1, \dots, x_{1728}$ .



**Figure 2.2:** Scatter plot of  $x_1, \dots, x_{1728}$ .

This picture gives an idea about completeness: in fact the step structure of the graph let us understand that reliable information about weakest earthquakes is more recent instead of that on strongest ones, which is to be found far back in time. Moreover we can also see the difference in return periods from seism with low magnitude and those with highest one.

Variable	Obs	Mean	Std. Dev.	Min	Max
MwDef	1728	4.701619	.5921769	3.95	7.32
Year	1728	1933.653	87.1905	1117	2014

**Figure 2.3:** Minimum and maximum value of magnitude and year of observation.

The following step is grouping observations in blocks of length one year. Since in many years data are not available, we obtain only 270 block maxima. Remember that independent records lead to independent block maxima. Now we can apply to  $z_1, \dots, z_{270}$  the maximum likelihood estimation method.

```
ML fit of GEV                                Number of obs =      270
Log likelihood = -232.94435                    Wald chi2(0) =      .
                                                Prob > chi2 =      .
```

blockmaxima	Coef.	Std. Err.	z	P> z	[95% Conf. Interval]
scale					
_cons	.5068625	.0253133	20.02	0.000	.4572493 .5564756
shape					
_cons	-.0722326	.0491349	-1.47	0.142	-.1685353 .02407
location					
_cons	5.381541	.0350767	153.42	0.000	5.312792 5.45029

Figure 2.4: Parameter estimates and 95% confidence intervals.

At first sight we are induced thinking that  $\hat{\xi} = -0.0722326 < 0$  means that the limit distribution of maxima can be represented by Weibull family, but the confidence interval extends well above zero. So, the evidence from data for a bounded distribution is not strong. Let's check the fitted model.

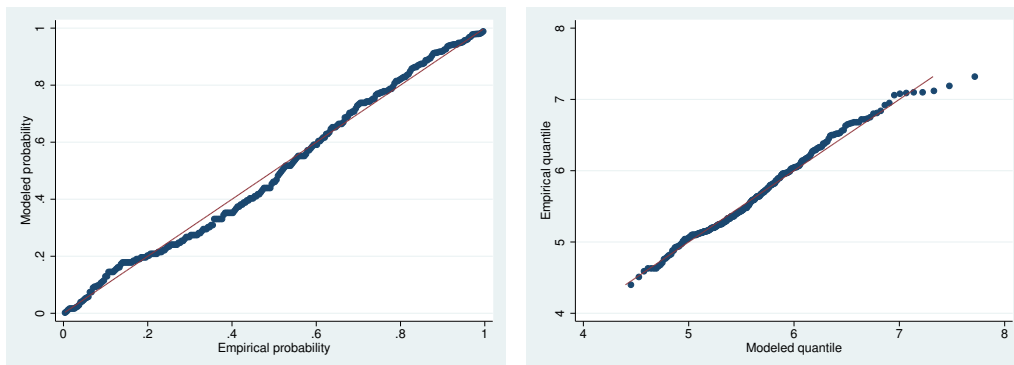
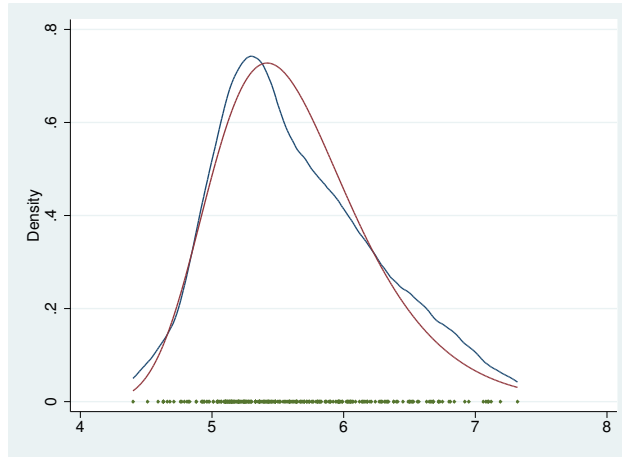


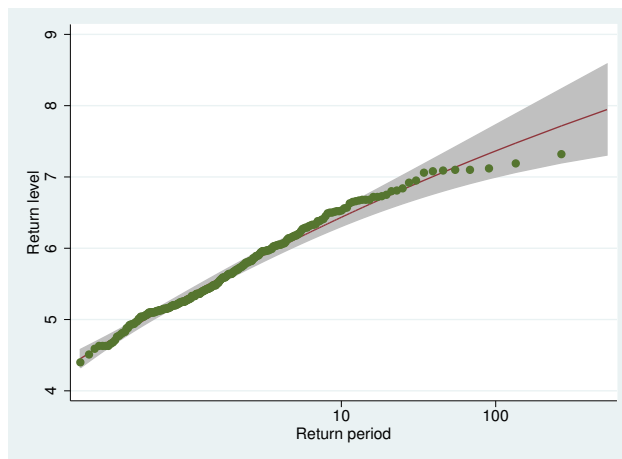
Figure 2.5: Probability plot (on the left) and Quantile plot (on the right).

Neither the probability plot nor the quantile plot give cause to doubt the validity of the fitted model: each set of plotted points is near linear. Even the comparison between kernel density (that is the "empirical" density obtained by an independent and identically distributed sample using a smoothing process) and fitted density doesn't instill doubt that model doesn't fit well.



**Figure 2.6:** Kernel density plot.

Finally we can use the return level plot for judging the effect of the estimated shape parameter. Since the model-based curve and empirical estimates are in reasonable agreement, we can thus deduce that the model is suitable.



**Figure 2.7:** Return level plot.



These diagnostic plots all suggest that model fits well but we are not still sure if the distribution family is the Weibull or the Gumbel one. So we resort to using the log-likelihood ratio test applied to the estimated model and a model with the constraint  $\xi = 0$ .

```

Likelihood-ratio test
(Assumption: gumbelitalia nested in gevitalia)
LR chi2(1) = 2.06
Prob > chi2 = 0.1507

Akaike's information criterion and Bayesian information criterion

```

Model	Obs	ll(null)	ll(model)	df	AIC	BIC
<u>gumbelitalia</u>	270	.	-233.9767	2	471.9534	479.1502
<u>gevitalia</u>	270	.	-232.9444	3	471.8887	482.684

Note: N=Obs used in calculating BIC; see [\[R\] BIC note](#)

**Figure 2.8:** Likelihood ratio test between the fitted model and the model representing the Gumbel family.

This test, which compares double difference between the two models log-likelihood to value of the  $\chi^2(1)$  statistic, results not significant. This means that fitted model doesn't produce more information than the simpler one, so the model with  $\xi = 0$  is better than the more general. Also AIC and BIC values confirm it.

This implies that block maxima distribution is well represented by the Gumbel family.

### 2.3.2 Ten years blocks

Let's now grouping the 1728 reliable observations in blocks of length 10 years: block maxima are now  $z_1, \dots, z_{52}$  and, as in the previous case, there are decades in which there aren't records.

From the following table we can see that parameter estimates are now all significant, including the shape parameter  $\hat{\xi} = -0.4136025 < 0$ . It's confidence interval doesn't include zero so we deduce that the block maxima distribution is now the Weibull distribution. Location and scale parameter estimates remain almost similar to the previous estimates.

```

ML fit of GEV
Log likelihood = -37.34228
Number of obs = 52
Wald chi2(0) = .
Prob > chi2 = .

```

blockmaxi~10	Coef.	Std. Err.	z	P> z	[95% Conf. Interval]
scale					
_cons	.53383	.0606337	8.80	0.000	.4149901 .65267
shape					
_cons	-.4136025	.0966682	-4.28	0.000	-.6030687 -.2241363
location					
_cons	6.19984	.081515	76.06	0.000	6.040073 6.359606

Figure 2.9: Parameter estimates for decade block maxima.

The model adequacy can be assessed using diagnostic plots. Both probability and quantile plots are roughly linear, suggesting good fit of the model.

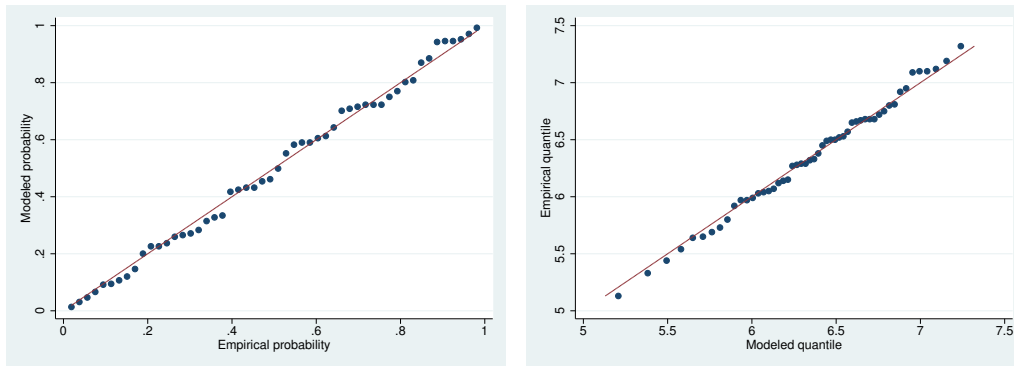
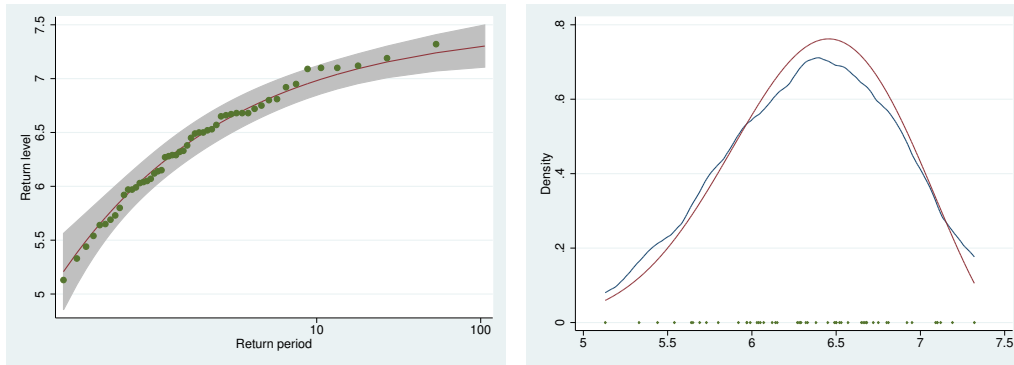


Figure 2.10: Probability plot (on the left) and Quantile plot (on the right) for block of length 10 years.

In the return level plot below we can see that the model-based curve and empirical estimates are in agreement, suggesting pertinence of the fitted model. Also the comparison between kernel and fitted density leads to the conclusion that model describes well the data.



**Figure 2.11:** Return level plot (on the left) and kernel density plot (on the right) for block of length 10 years.

### 2.3.3 Fifty years blocks

Blocks are now of 50 years length, so the number of records we use to perform our analysis is very low leading to probable large estimation variance. Let them be  $z_1, \dots, z_{14}$ .

```

ML fit of GEV                                Number of obs =      14
Log likelihood = -3.4930601                    Wald chi2(0) =      .
                                                Prob > chi2 =      .

```

blockmaxi~50	Coef.	Std. Err.	z	P> z	[95% Conf. Interval]
scale					
_cons	.385448	.0976931	3.95	0.000	.193973 .5769231
shape					
_cons	-.6642924	.2265895	-2.93	0.003	-1.1084 -.2201852
location					
_cons	6.772366	.1129873	59.94	0.000	6.550915 6.993817

**Figure 2.12:** Parameter estimates for fifty years block maxima.

As in the ten years blocks case, parameter estimates are all significant. In particular we have  $\hat{\xi} = -0.6642924 < 0$  and it's confidence interval doesn't contain zero leading to the Weibull distribution family.

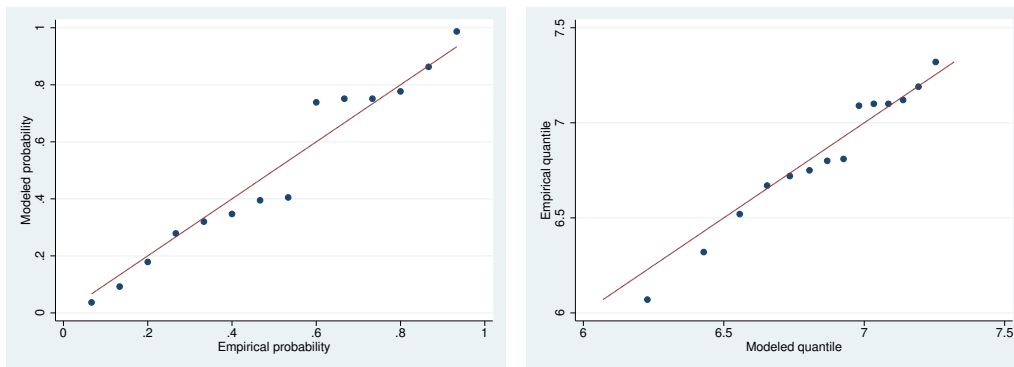
Conversely from what we expected, the variance of estimates is not too large, as we can see in the table below.

Covariance matrix of coefficients of gevfit model

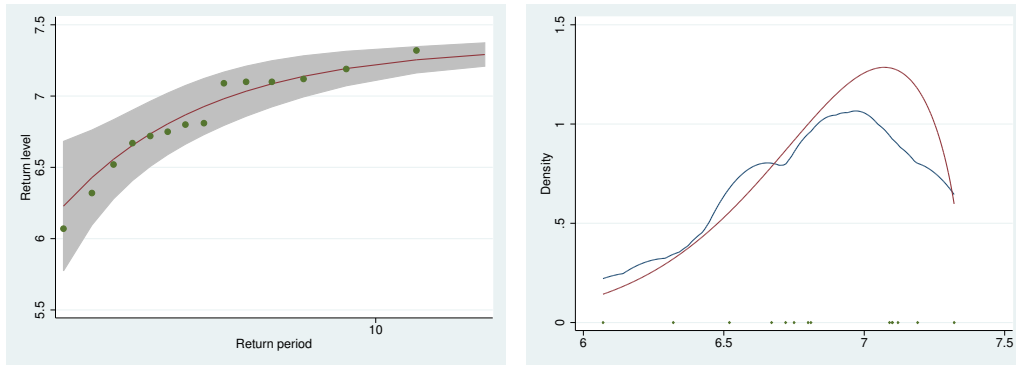
	e (V)	scale _cons	shape _cons	location _cons
scale _cons		.00954395		
shape _cons		-.01462085	.05134278	
location _cons		-.00394448	-.01046023	.01276614

**Figure 2.13:** Variance and covariance matrix of parameter estimates for 50 years block maxima.

Adequacy of fitted model can be assessed by diagnostic plots. Although we have only 14 observations, probability and quantile plots are still roughly linear. In the return level plot empirical estimates adapt quite well to the model based-curve and the kernel density shape is vaguely similar to the fitted model density.



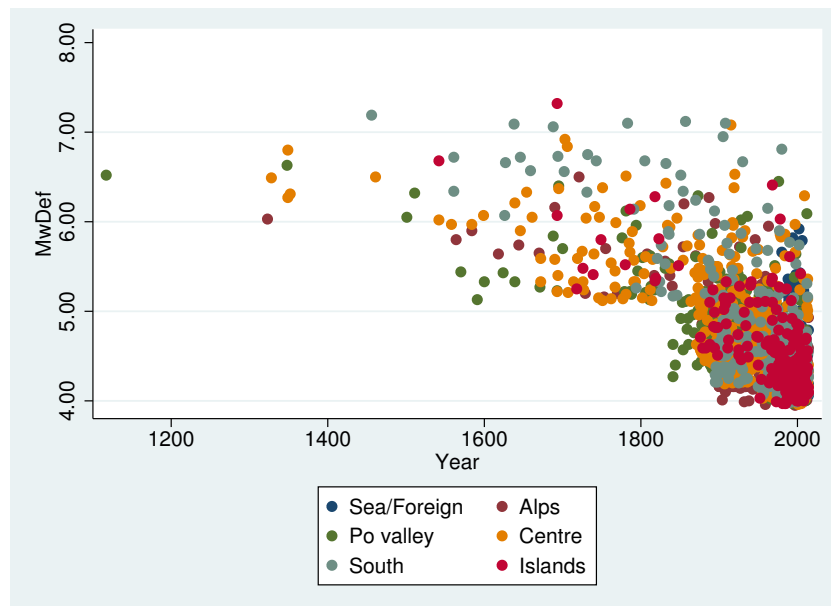
**Figure 2.14:** Probability plot (on the left) and Quantile plot (on the right) for block of length 50 years.



**Figure 2.15:** Return level plot (on the left) and kernel density plot (on the right) for block of length 50 years.

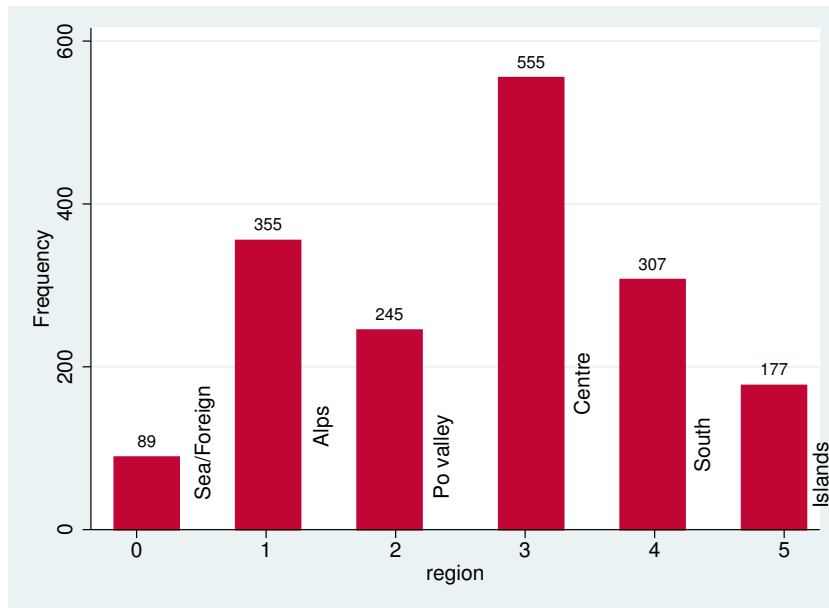
### 2.3.4 Zone analysis

In previous subsections we studied how parameter estimates change modifying the block size. Now our aim is understanding what is the family distribution of regional one year block maxima.



**Figure 2.16:** Scatter plot of regional records.

From an introductory investigation we see that number of records substantially vary among regions: this is due not only for a difference in seismic, but mainly from data reliability which changes region to region as explained at the beginning of the chapter. For example the first observation in the "Sea/Foreign" region dates back to 1980, instead the "Po valley" one dates back to 1117 (it's the first record of the whole cleaned dataset).



**Figure 2.17:** Frequency of regional records.

Let's take a look to the following tables.

Regional scale and shape parameter estimates are all significant and values are similar moving from a zone to another: in fact  $\hat{\mu} \in [4.586424, 5.000841]$  and  $\hat{\sigma} \in [0.3923321, 0.5900632]$ .

Different is the behaviour of the shape parameter estimates, which are little significant or totally insignificant in each region considered.

For "Sea/Foreign" region we have  $\hat{\xi} = -0.3041806 < 0$  and it's confidence interval doesn't contain zero so, despite the parameter estimate is little significant, we can deduce that block maxima distribution is well represented by the Weibull family.

For "Alps" region  $\hat{\xi} = 0.025207 > 0$  but it isn't significant at all and it's

confidence interval extends well below zero. Fitting the model with the constraint  $\xi = 0$ , we see that  $\hat{\mu}$  and  $\hat{\sigma}$  assume values similar to those in the more general model and are still significant, thus we can deduce that a Gumbel family is more appropriate.

```

-> region = Sea/Foreign

ML fit of GEV                                Number of obs =      29
Log likelihood = -13.727656                    Wald chi2(0) =      .
                                                Prob > chi2 =      .

```

blockmaxim~g	Coef.	Std. Err.	z	P> z	[95% Conf. Interval]
scale					
_cons	.3923321	.0587779	6.67	0.000	.2771295 .5075346
shape					
_cons	-.3041806	.1394031	-2.18	0.029	-.5774056 -.0309557
location					
_cons	5.000841	.0815326	61.34	0.000	4.84104 5.160642

---

```

-> region = Alps

ML fit of GEV                                Number of obs =     144
Log likelihood = -104.38939                    Wald chi2(0) =      .
                                                Prob > chi2 =      .

```

blockmaxim~g	Coef.	Std. Err.	z	P> z	[95% Conf. Interval]
scale					
_cons	.4179327	.0308849	13.53	0.000	.3573994 .4784661
shape					
_cons	.025207	.0801805	0.31	0.753	-.1319439 .1823579
location					
_cons	4.586424	.0409042	112.13	0.000	4.506254 4.666595

**Figure 2.18:** Regional parameter estimates and 95% confidence intervals.

Instead "Po valley" region has a negative shape parameter estimate  $\hat{\xi} = -0.0942388 < 0$ , it discloses the same situation of "Alps" region, with a confidence interval containing zero and high p-value. Performing a model with  $\xi = 0$ ,  $\hat{\mu}$  and  $\hat{\sigma}$  values remain very similar inducing us choosing also in this case a Gumbel family distribution than the Fréchet one.

Regarding "Centre" region we can see a little significance of  $\hat{\xi} = -0.1371585 < 0$ , but zero doesn't belong to the confidence interval, than we can say that in this case block maxima follow the Weibull distribution.

```

-> region = Po valley

ML fit of GEV                                Number of obs =      143
Log likelihood = -112.92514                    Wald chi2(0) =      .
                                                Prob > chi2 =      .

```

blockmaxim~g	Coef.	Std. Err.	z	P> z	[95% Conf. Interval]
scale					
_cons	.4791222	.031988	14.98	0.000	.4164268 .5418176
shape					
_cons	-.0942388	.0605681	-1.56	0.120	-.21295 .0244725
location					
_cons	4.784946	.0450312	106.26	0.000	4.696687 4.873206

```

-> region = Centre

ML fit of GEV                                Number of obs =      195
Log likelihood = -186.32241                    Wald chi2(0) =      .
                                                Prob > chi2 =      .

```

blockmaxim~g	Coef.	Std. Err.	z	P> z	[95% Conf. Interval]
scale					
_cons	.5771983	.0342522	16.85	0.000	.5100652 .6443314
shape					
_cons	-.1371585	.0585835	-2.34	0.019	-.25198 -.022337
location					
_cons	4.966685	.0471752	105.28	0.000	4.874223 5.059146

**Figure 2.19:** Regional parameter estimates and 95% confidence intervals.

Finally "South" and "Islands" regions show the same shape parameter estimates behaviour: in both cases we have a not significant and positive value of  $\hat{\xi}$  (precisely  $\hat{\xi} = 0.1389487$  for "South",  $\hat{\xi} = 0.1092782$  for "Islands") with confidence intervals extending below zero. As in previous cases, scale and location parameter estimates of restricted model don't differ so much from estimates of the more general one and are still significant, making us



prone to choose the Gumbel family distribution for block maxima.

---

```

-> region = South

ML fit of GEV                                Number of obs =      134
Log likelihood = -152.02562                    Wald chi2(0) =      .
                                                Prob > chi2 =      .

```

blockmaxim~g	Coef.	Std. Err.	z	P> z	[95% Conf. Interval]	
scale						
_cons	.5900632	.0478276	12.34	0.000	.4963228	.6838035
shape						
_cons	.1389487	.0906265	1.53	0.125	-.0386759	.3165733
location						
_cons	4.844757	.0606206	79.92	0.000	4.725943	4.963571

---

```

-> region = Islands

ML fit of GEV                                Number of obs =      91
Log likelihood = -64.841048                    Wald chi2(0) =      .
                                                Prob > chi2 =      .

```

blockmaxim~g	Coef.	Std. Err.	z	P> z	[95% Conf. Interval]	
scale						
_cons	.3948911	.0357114	11.06	0.000	.3248982	.4648841
shape						
_cons	.1092782	.0838482	1.30	0.192	-.0550613	.2736176
location						
_cons	4.643468	.0469782	98.84	0.000	4.551393	4.735544

Figure 2.20: Regional parameter estimates and 95% confidence intervals.

### 2.3.5 Results comparison

In order to compare in a simple manner results obtained in previous subsections, it can be useful summarize them in a table<sup>2</sup>.

---

<sup>2</sup>one star (\*) if  $p < 0.05$ , two stars (\*\*) if  $p < 0.01$ , and three stars (\*\*\*) if  $p < 0.001$ .

Zone	shape $\hat{\xi}$	95% C.I.	Dist.Fam.
Italy (1 year b.m.)	-0.0722326	[-0.1685353 , 0.02407]	Gumbel
Italy (10 year b.m.)	-0.4136025 ***	[-0.6030687 , -0.2241363]	Weibull
Italy (50 year b.m.)	-0.6642924 **	[-1.1084 , -0.2201852]	Weibull
"Sea/Foreign" region	-0.3041806 *	[-0.5774056 , -0.0309557]	Weibull
"Alps" region	0.025207	[-0.1319439 , 0.1823579]	Gumbel
"Po valley" region	-0.0942388	[-0.21295 , 0.0244725]	Gumbel
"Centre" region	-0.1371585 *	[-0.25198 , -0.022337]	Weibull
"South" region	0.1389487	[-0.0386759 , 0.3165733]	Gumbel
"Islands" region	0.1092782	[-0.0550613 , 0.2736176]	Gumbel

## 2.4 GPD parameters estimation

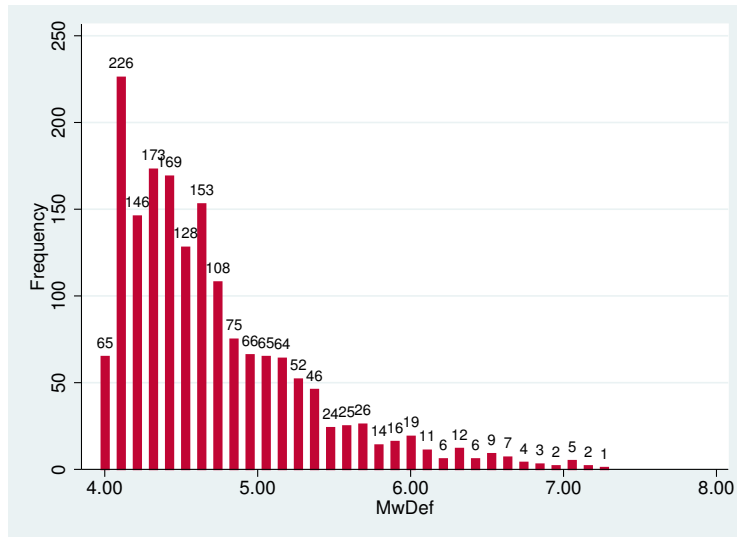
### 2.4.1 Global analysis

Let us consider the iid sample formed by the 1728 records of the cleaned dataset: remember that those data refer to whole Italian territory.

The first step consists in ordering these values in an increasing manner respect to the magnitude  $x_{(1)} < \dots < x_{(1728)}$  and defining threshold excess as

$$y_i = x_{(i)} - u, \quad i = 1, \dots, 1728$$

for a certain threshold  $u$ .



**Figure 2.21:** Histogram of Italian records for different values of magnitude.

Default options put threshold at the minimum recorded value of magnitude, corresponding to  $u = 3.95$ , leading to these parameter estimates:

Taking sample minimum of 3.95 as the threshold.

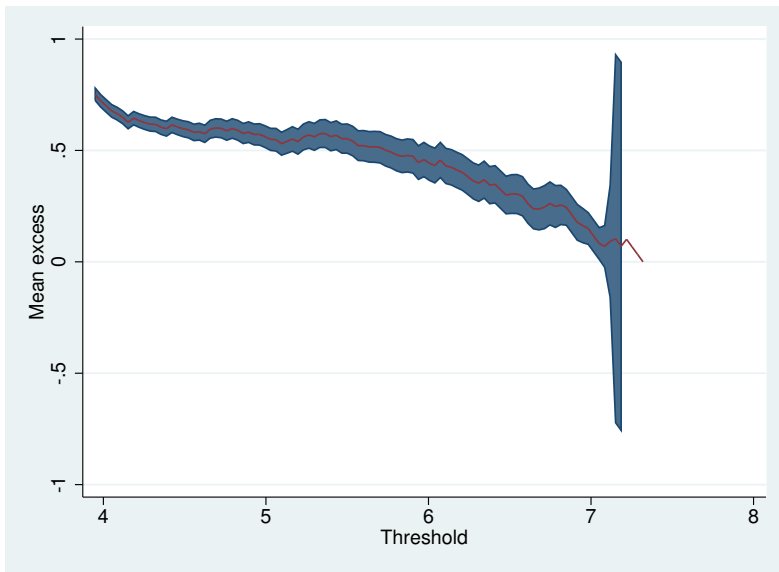
```
ML fit of generalized Pareto distribution      Number of obs =      1727
Log likelihood = -1189.1326                  Wald chi2(0) =      .
                                              Prob > chi2 =      .
```

MwDef	Coef.	Std. Err.	z	P> z	[95% Conf. Interval]
lnsig					
_cons	-.0897912	.0294933	-3.04	0.002	-.147597 - .0319853
xi					
_cons	-.2216549	.017868	-12.41	0.000	-.2566756 - .1866342
sig	.9141221	.0269605			.8627787 .9685208

**Figure 2.22:** Parameter estimates and 95% confidence intervals for  $u = 3.95$ .

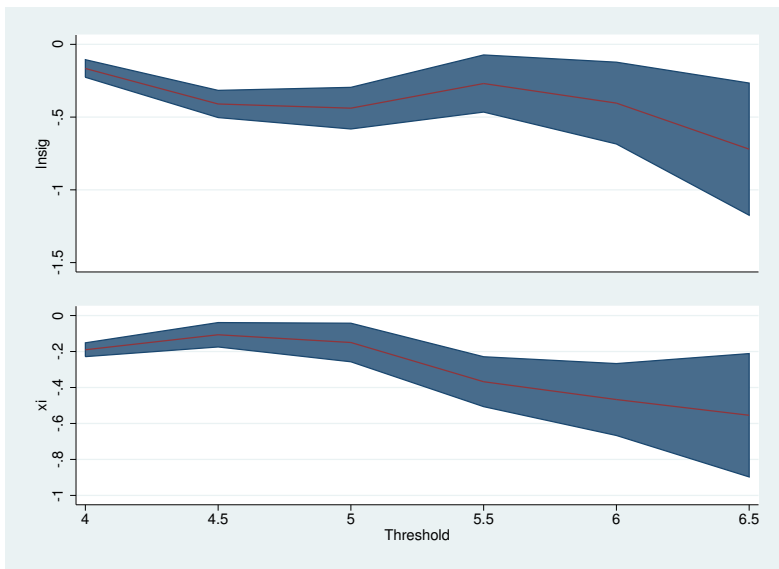
Location parameter estimate doesn't appear because of the assumption  $\tilde{\mu} = 0$  made in Section 1.6. Scale and shape parameter estimates are both significant, in particular  $\hat{\xi} = -0.2216549 < 0$  with a 95% confidence interval not including zero.

From the histogram above, we could think to choose the initial value of threshold as  $u = 4.5$ . In fact this choice would leave almost half observations below and half above  $u$ , avoiding violation of model asymptotic bases (if the threshold is too low) and high variance (if the threshold is too high). More accurate instrument in threshold selection, which uses threshold excesses  $y_i$  for  $i = 1, \dots, 1728$ , is the mean residual life plot: above the threshold at which the generalised Pareto distribution provides a valid approximation to the excess distribution function  $F_u(x)$ , the mean residual life plot should be approximately linear (see properties of GPD in Section 1.6). For  $u < 4.5$  the graph below appears to curve, for  $4.5 \leq u \leq 5$  is approximately linear, and for  $u > 5$  it decreases rapidly. We can be tempted to choose as threshold  $u = 5$ , but above this value there are only 419 records. Thus is better to choose as threshold  $u = 4.5$ .



**Figure 2.23:** Mean residual life plot.

Support for this choice is provided by the following graph, representing parameter estimates at different values of the threshold. Perturbations for high thresholds are visible but they are small relative to sampling errors, so  $u = 4.5$  appears reasonable.



**Figure 2.24:** Parameter estimates against threshold for Italian data.

Maximum likelihood estimates in this case are given in the table below.

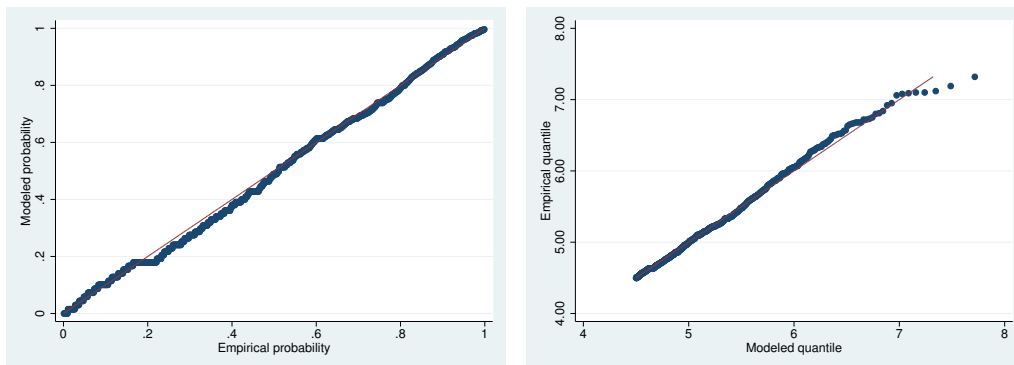
```
ML fit of generalized Pareto distribution      Number of obs =      918
Log likelihood = -443.80699                  Wald chi2(0) =      .
                                              Prob > chi2 =      .
```

MwDef	Coef.	Std. Err.	z	P> z	[95% Conf. Interval]
lnsig					
_cons	-.4097495	.0476707	-8.60	0.000	-.5031825 - .3163166
xi					
_cons	-.1068006	.0345775	-3.09	0.002	-.1745711 - .03903
sig	.6638165	.0316446			.6046035 .7288287

**Figure 2.25:** Parameter estimates and 95% confidence intervals for  $u = 4.5$ .

Parameter estimates are both significant and the confidence interval for  $\hat{\xi} = -0.1068006 < 0$  is in the negative domain, leading to the Weibull distribution family.

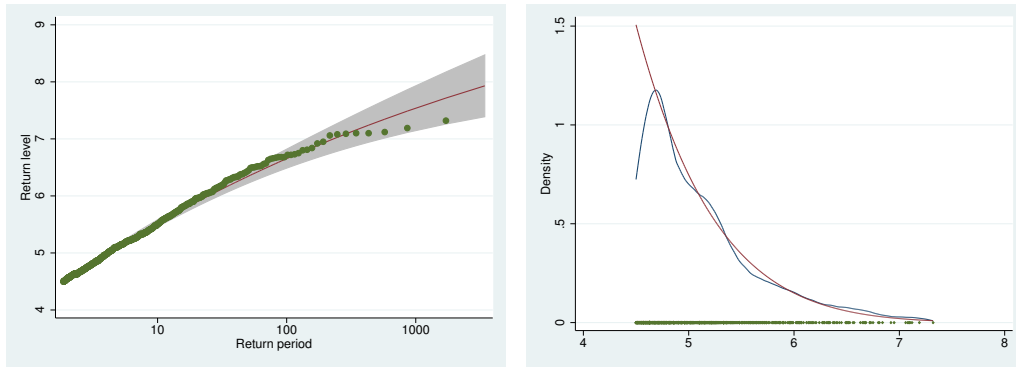
Same diagnostic plots of GEV distribution case can be used for assessing model adequacy.



**Figure 2.26:** Probability plot (on the left) and Quantile plot (on the right) for GPD with  $u = 4.5$ .

As we can see, probability and quantile plot are pretty much linear, indicating model good fit. The return level plot shows that the model-based curve

and empirical estimates are in reasonable agreement and the kernel density plot, even if the empirical density is initially concave, leaves no doubt on model adequacy.

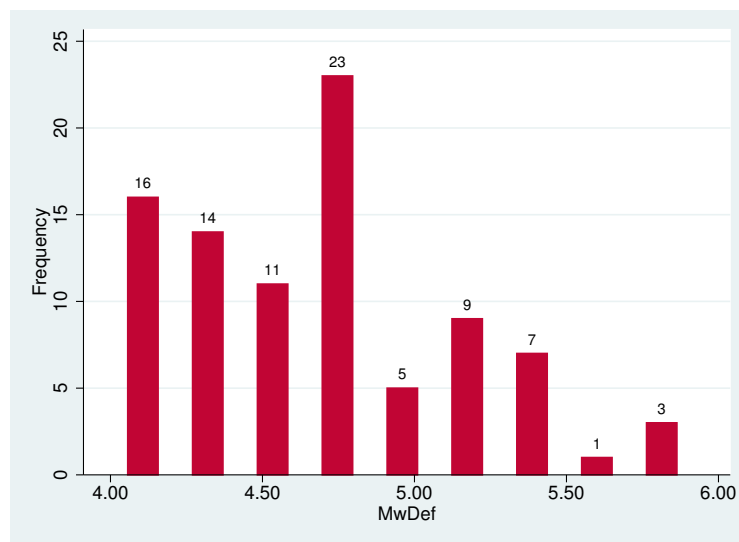


**Figure 2.27:** Return level plot (on the left) and kernel density plot (on the right) for GPD with  $u = 4.5$ .

## 2.4.2 Zone analysis

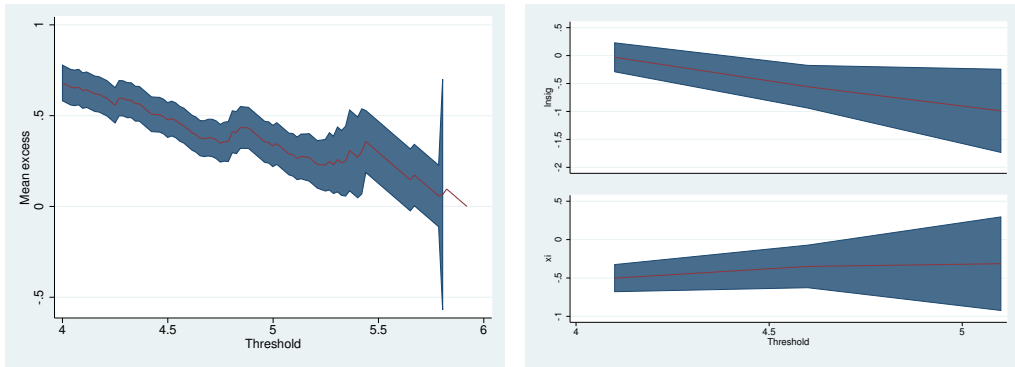
As for the GEV distribution, we want to perform a differentiating analysis by zones also for the generalised Pareto distribution.

The first zone we analyse is "Sea/Foreign" region for which we have 89 observations distributed as follow:



**Figure 2.28:** Histogram of "Sea/Foreign" region records.

Minimum and maximum magnitude recorded are 4 and 5.92 respectively and the histogram suggests to choose the threshold at 4.6. This value in fact allows to leave almost half observations below and half above  $u$ . Mean residual life plot is approximately linear for  $u \in [4.5, 4.75]$  and plotting parameter estimates against threshold we can see that sampling errors increase beyond the value 4.6; thus these graphs give support to our choice.



**Figure 2.29:** Mean residual life plot (on the left) and parameter estimates against threshold (on the right) for "Sea/Foreign" region.

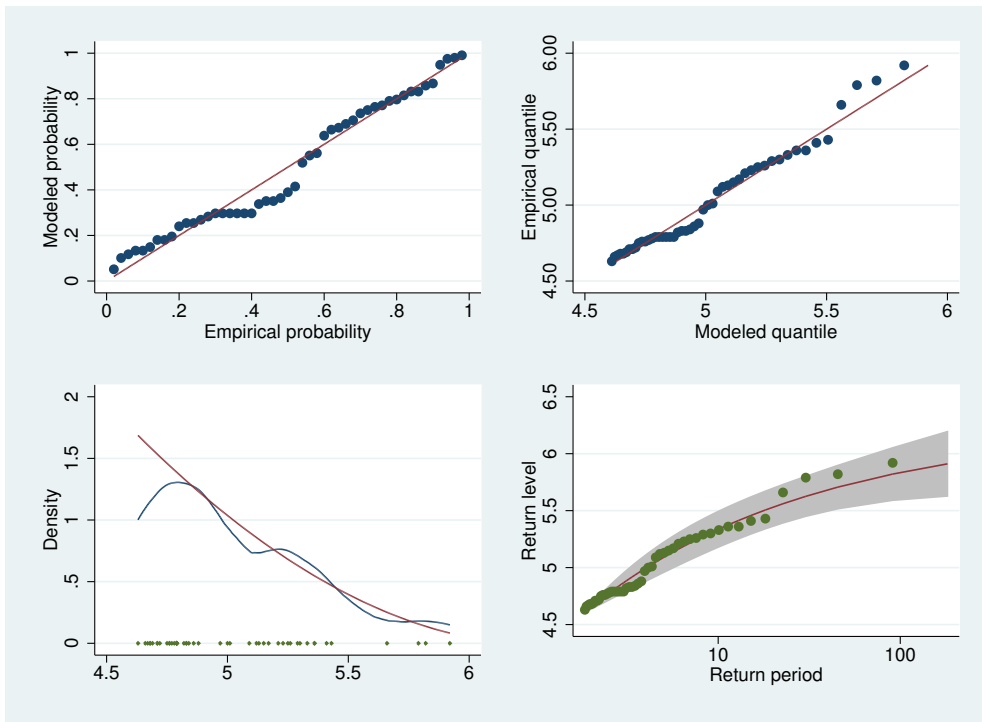
We can now estimate scale and shape parameter for the threshold selected:

```
ML fit of generalized Pareto distribution           Number of obs   =           49
Log likelihood = -4.5695061                       Wald chi2(0)    =           .
                                                    Prob > chi2     =           .
```

MwDef	Coef.	Std. Err.	z	P> z	[95% Conf. Interval]	
lnsig	_cons	-.5570596	.194759	-2.86	0.004	-.9387802 - .175339
	xi	-.3496852	.1414864	-2.47	0.013	-.6269934 - .072377
	sig	.5728911	.1115757			.3911046 .8391725

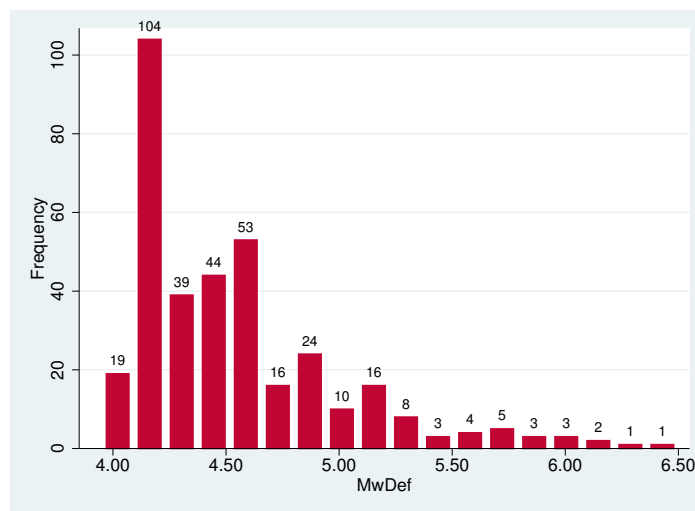
**Figure 2.30:** Parameter estimates and 95% confidence intervals for  $u = 4.6$ .

Shape parameter estimate  $\hat{\xi} = -0.3496852 < 0$  isn't much significant, however it's confidence interval doesn't contain zero, leading to the Weibull distribution family. All model diagnostic plots, even if with some discrepancies, confirm model adequacy.



**Figure 2.31:** Probability plot and Quantile plot (on the top), Kernel density plot and Return level plot (on the bottom) for  $u = 4.6$ .

For "Alps" region we have 355 observations.

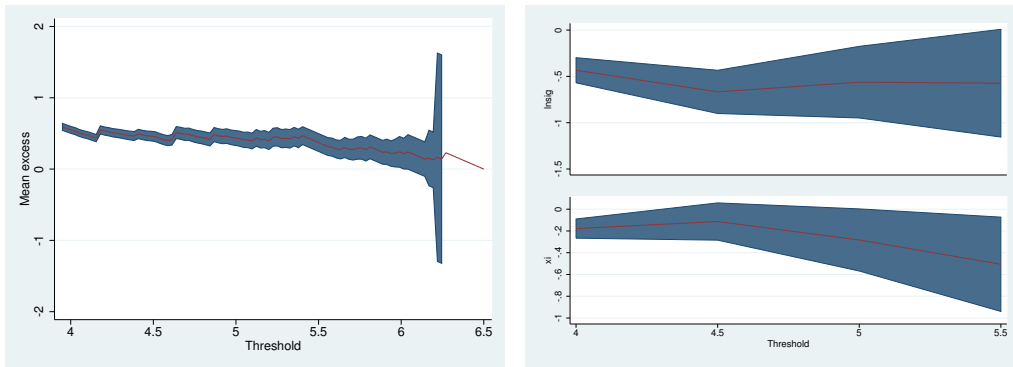


**Figure 2.32:** Histogram of "Alps" region records.



The histogram above shows how they are distributed. Minimum magnitude is 3.95 and the maximum is 6.5. Furthermore it suggests to choose as threshold  $u = 4.4$ .

The mean residual life plot is approximately linear for  $u \in [4.3, 4.6]$ : this isn't the unique interval but the others concern too high values of threshold, with scarce number of excesses which can leads to high variance. In fact sampling errors in the graph of parameter estimates against threshold increase up to  $u = 4.5$ .



**Figure 2.33:** Mean residual life plot (on the left) and parameter estimates against threshold (on the right) for "Alps" region.

Thus we can choose the threshold level at  $u = 4.4$  and give maximum likelihood parameter estimates:

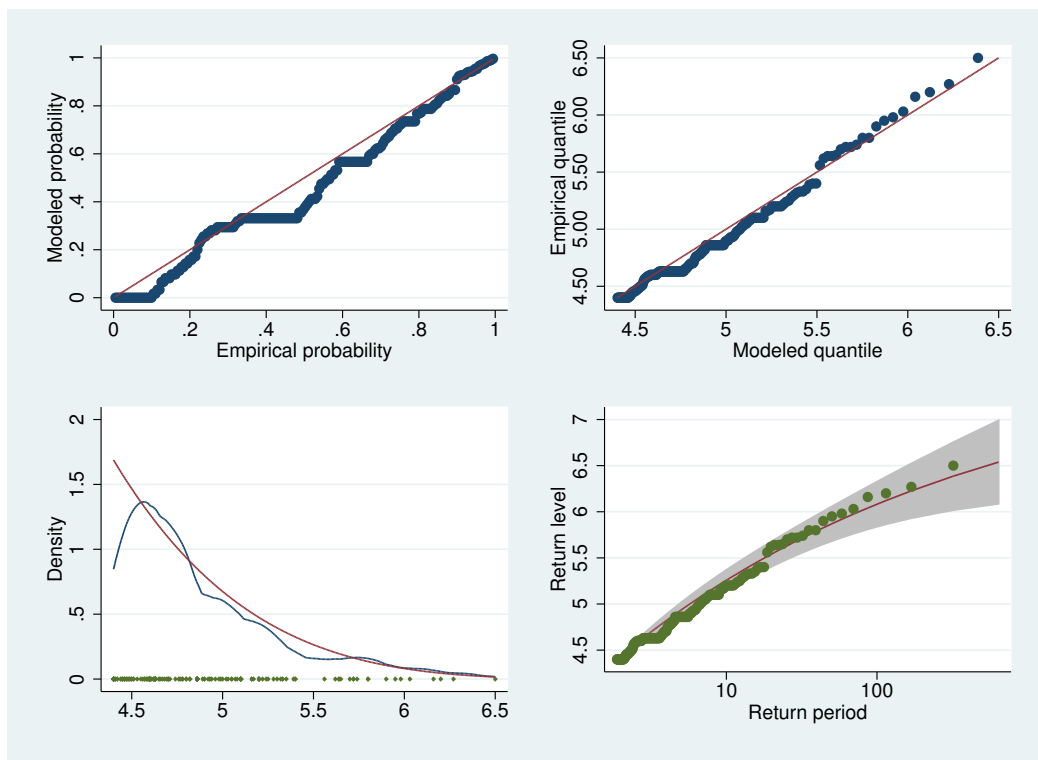
```
ML fit of generalized Pareto distribution           Number of obs =          171
Log likelihood = -50.622055                       Wald chi2(0) =           .
                                                    Prob > chi2 =            .
```

MwDef	Coef.	Std. Err.	z	P> z	[95% Conf. Interval]
lensig					
_cons	-.5234722	.1036739	-5.05	0.000	-.7266692    -.3202751
xi					
_cons	-.1804924	.07135	-2.53	0.011	-.3203359    -.0406489
sig	.5924599	.0614226			.4835168    .7259493

**Figure 2.34:** Parameter estimates and 95% confidence intervals for  $u = 4.4$ .

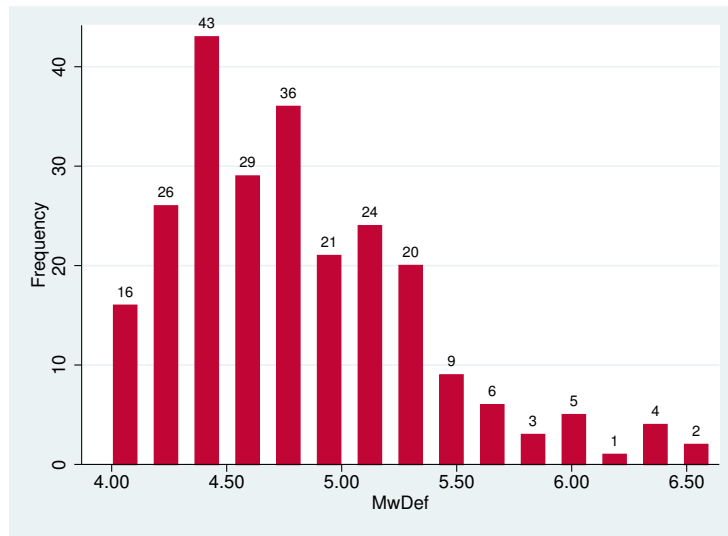
From the table we can see that scale parameter estimate is highly significant, unlike the shape one which is little significant. However  $\hat{\xi} = -0.1804924$  is negative and has a confidence interval completely contained in the negative domain so, also in this case, the family distribution is the Weibull.

As usually we use diagnostic plots for assessing model adequacy, which is pretty good even if the goodness-of-fit in the probability plot seems unconvincing.



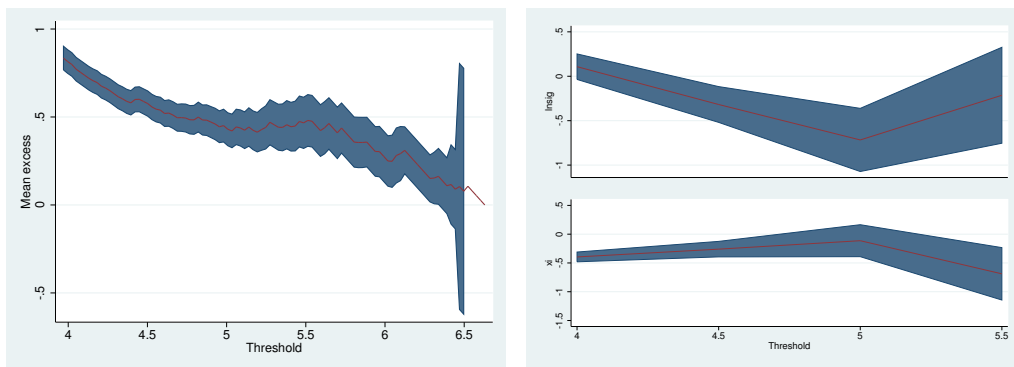
**Figure 2.35:** Probability plot (on the top left), Quantile plot (on the top right), Kernel density plot (on the bottom left) and Return level plot (on the bottom right) for  $u = 4.4$ .

For "Po valley" region we have 245 records, distributed as the histogram in the next page shows. At first sight a good threshold is  $u = 4.7$ .



**Figure 2.36:** Histogram of "Po valley" region records.

Mean residual life plot and the other graph below both suggest choosing as threshold a value  $u \in [4.5, 4.6]$ .



**Figure 2.37:** Mean residual life plot (on the left) and parameter estimates against threshold (on the right) for "Po valley" region.

So maximum likelihood estimates for  $u = 4.6$  are the following:

ML fit of generalized Pareto distribution		Number of obs	=	146		
Log likelihood = -50.406237		Wald chi2(0)	=	.		
		Prob > chi2	=	.		

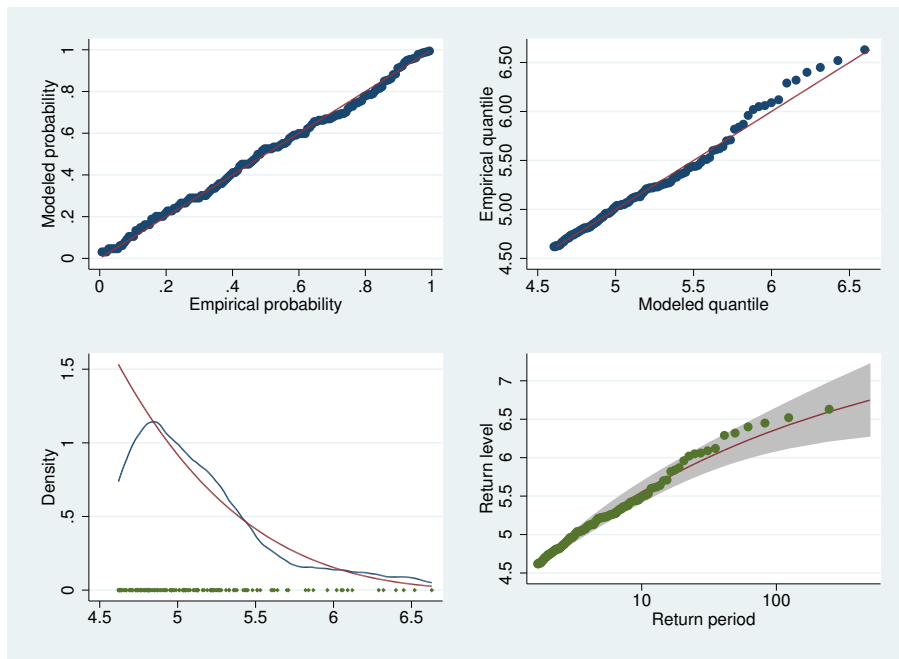
  

MwDef	Coef.	Std. Err.	z	P> z	[95% Conf. Interval]	
lnsig						
_cons	-.4522351	.1134558	-3.99	0.000	-.6746044	-.2298658
xi						
_cons	-.2025167	.0793966	-2.55	0.011	-.3581311	-.0469023
sig	.6362046	.0721811			.5093579	.7946402

**Figure 2.38:** Parameter estimates and 95% confidence intervals for  $u = 4.6$ .

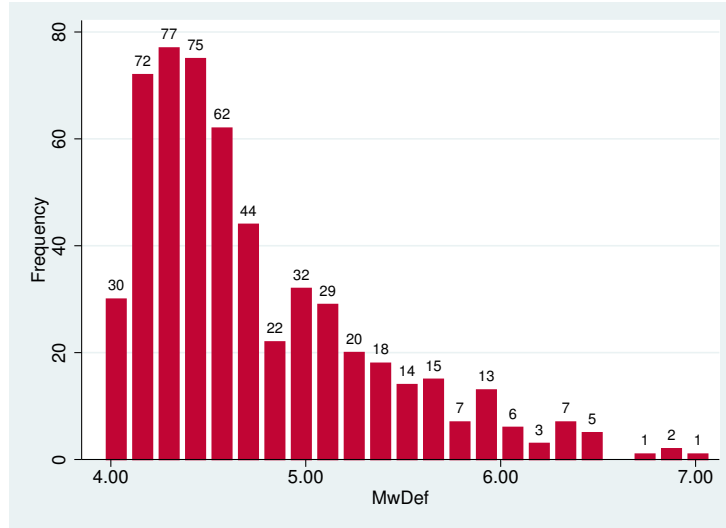
Shape parameter isn't much significant but it's confidence interval doesn't contain zero so, being  $\hat{\xi} = -0.2025167 < 0$ , we can say that the family distribution is again the Weibull one.

Even in this case model diagnostic plots confirm a good fit of the model.



**Figure 2.39:** Probability plot and Quantile plot (on the top), Kernel density plot and Return level plot (on the bottom) for  $u = 4.6$ .

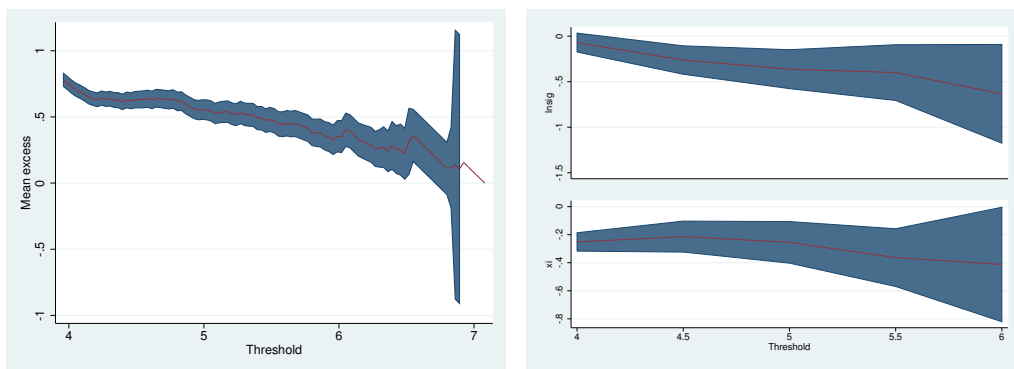
For "Centre" region there are 555 records distributed as follow:



**Figure 2.40:** Histogram of "Centre" region records.

The value of threshold which leaves half observations on the left and half on the right is approximately  $u = 4.5$ .

Using the mean residual life plot we can see that a good interval in which selecting the threshold could be  $[4.3, 4.8]$ , because there the graph is almost linear. Furthermore, plotting parameter estimates against threshold, we can see certain stability in estimates for  $u \in [4.5, 5]$ . Since sampling errors increase with threshold increasing, we can opt for  $u = 4.5$ .



**Figure 2.41:** Mean residual life plot (on the left) and parameter estimates against threshold (on the right) for "Centre" region.

Both parameter estimates for the selected threshold are significant: in

```

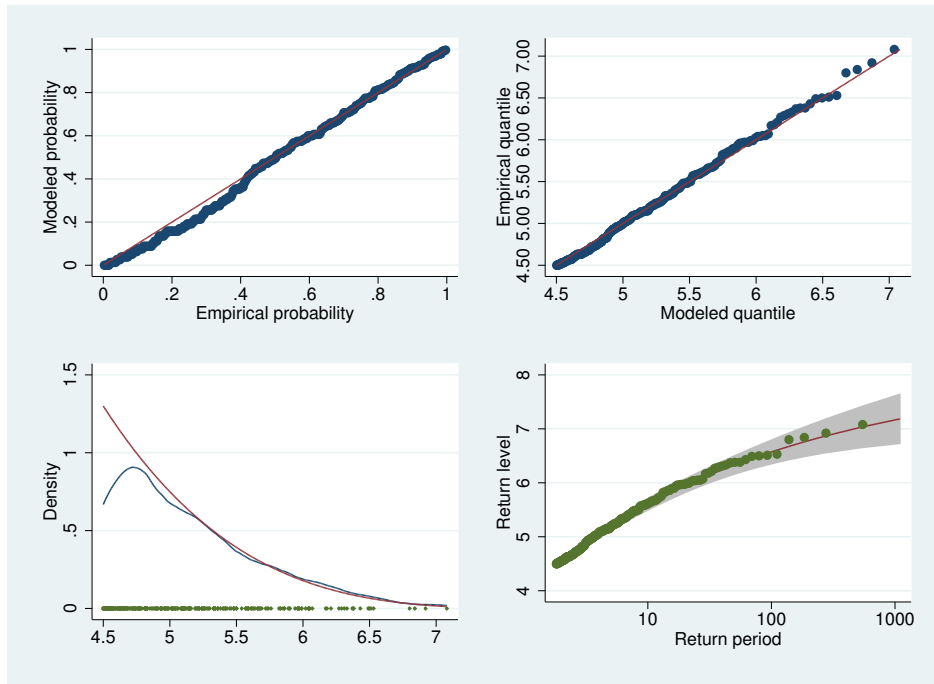
ML fit of generalized Pareto distribution          Number of obs =      301
Log likelihood = -157.73659                     Wald chi2(0) =      .
                                                Prob > chi2 =      .

```

MwDef	Coef.	Std. Err.	z	P> z	[95% Conf. Interval]	
lnsig						
_cons	-.2623753	.0796314	-3.29	0.001	-.4184499	-.1063006
xi						
_cons	-0.2135829	.0563068	-3.79	0.000	-.3239422	-.1032236
sig	.7692223	.0612543			.6580661	.8991543

**Figure 2.42:** Parameter estimates and 95% confidence intervals for  $u = 4.5$ .

particular we have  $\hat{\xi} = -0.2135829 < 0$  with a confidence interval not containing zero. Thus the family distribution is the Weibull.

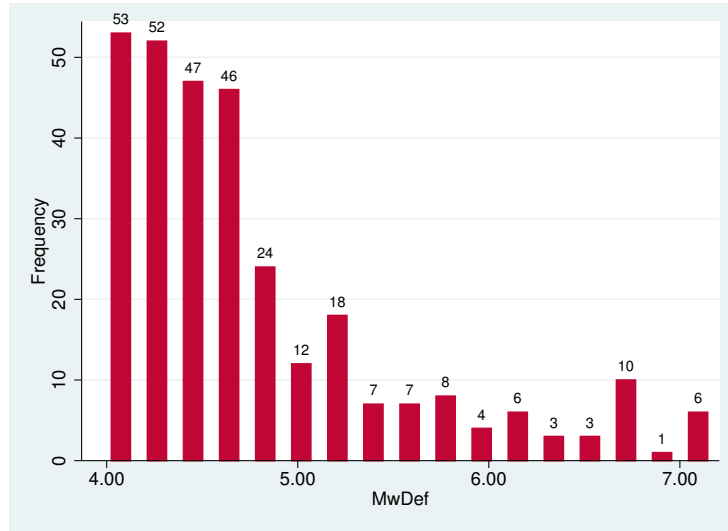


**Figure 2.43:** Probability plot and Quantile plot (on the top), Kernel density plot and Return level plot (on the bottom) for  $u = 4.5$ .

In every diagnostic plot we can see good agreement between empirical

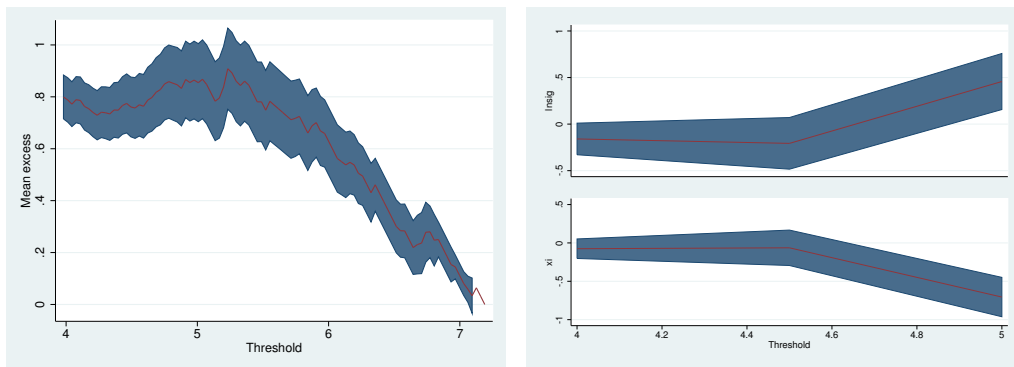
and fitted data which allows us to say that chosen model is adequate.

Let's now analyse "South" region, whose observations amount to 307.



**Figure 2.44:** Histogram of "South" region records.

From the first graph below we can see that the threshold  $u = 4.5$  suggested by the histogram is not good because mean residual life plot is not linear in a neighbourhood of this value, but it is approximately linear for  $u \in [4.6, 4.9]$ .



**Figure 2.45:** Mean residual life plot (on the left) and parameter estimates against threshold (on the right) for "South" region.

The second graph shows constant sampling errors for  $u > 4.5$ , then we could choose  $u = 4.6$  as threshold.

Remember that too high values of  $u$  can't be chosen because of an insufficient number of excesses.

For  $u = 4.6$  maximum likelihood estimates are both not significant and  $\hat{\xi} = -0.1068119 < 0$  has a confidence interval containing zero as the following table shows.

```
ML fit of generalized Pareto distribution      Number of obs   =      141
Log likelihood = -103.80488                  Wald chi2(0)    =      .
                                           Prob > chi2     =      .
```

MwDef	Coef.	Std. Err.	z	P> z	[95% Conf. Interval]
lnsig					
_cons	-.1569833	.1589475	-0.99	0.323	-.4685147 .154548
xi					
_cons	-0.1068119	.1351037	-0.79	0.429	-.3716103 .1579866
sig	.8547183	.1358553			.6259313 1.16713

Figure 2.46: Parameter estimates and 95% confidence intervals for  $u = 4.6$ .

In this situation we can think that a model with  $\xi = 0$  would be better: it isn't so. In fact in this case scale parameter estimate has p-value equal to one, that is estimate not significant at all, and also AIC and BIC values suggest the more general model would be better.

```
Likelihood-ratio test
(Assumption: gpdgumbelsouth nested in gpdsouth)
LR chi2(1) = 74.39
Prob > chi2 = 0.0000
```

```
Akaike's information criterion and Bayesian information criterion
```

Model	Obs	ll(null)	ll(model)	df	AIC	BIC
<u>gpdgumbels~h</u>	141	.	-141	1	284	286.9488
<u>gpdsouth</u>	141	.	-103.8049	2	211.6098	217.5073

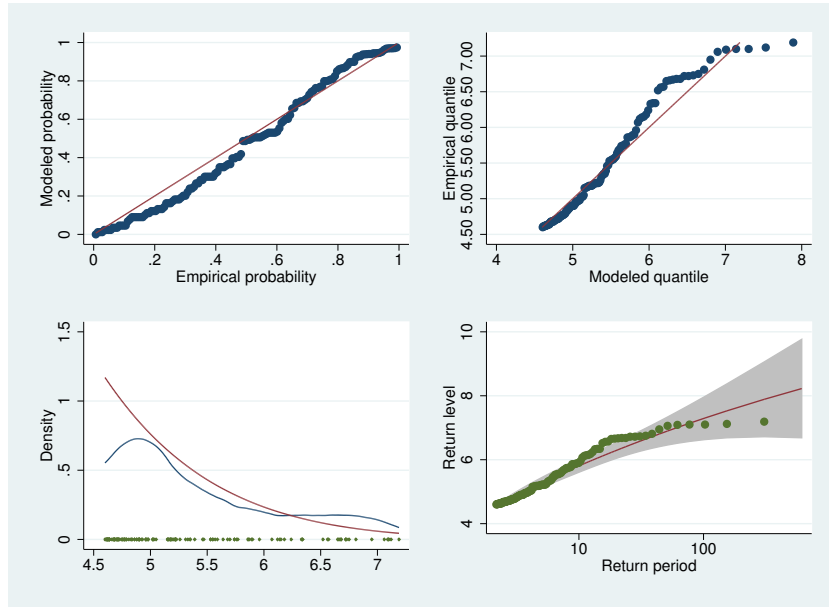
Note: N=Obs used in calculating BIC; see [R] BIC note

Figure 2.47: Likelihood ratio test between the fitted model and the model with  $\xi = 0$ .

This lack of fit can be seen also from the probability and quantile plot, in which only few points lie on the diagonal line, or from discrepancies between

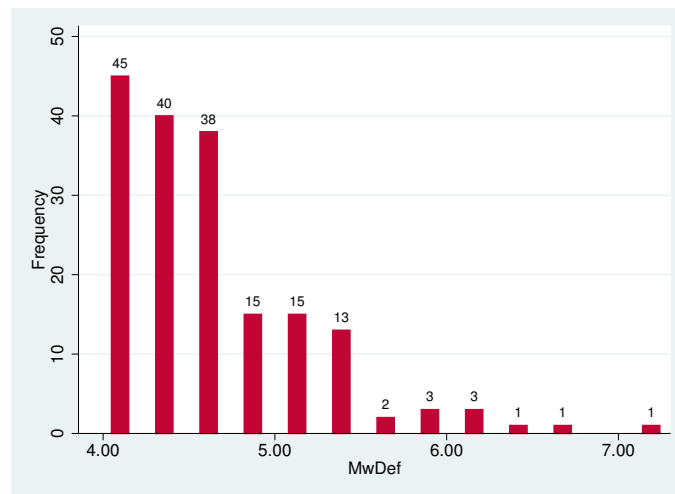


fitted and kernel density. In the return level plot observations lie even out of confidence interval.



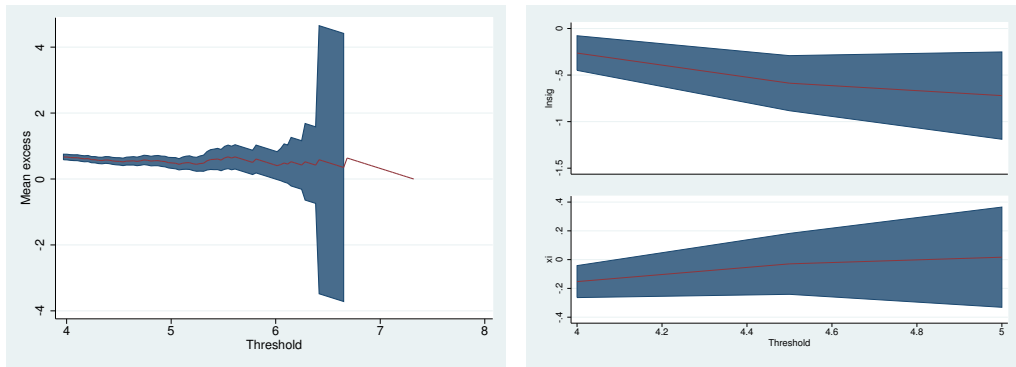
**Figure 2.48:** Probability plot (on the top left) and Quantile plot (on the top right), Kernel density plot (on the bottom left) and Return level plot (on the bottom right) for  $u = 4.6$ .

Finally the last region is "Islands", with 177 records.



**Figure 2.49:** Histogram of "Islands" region records.

Suggested threshold is  $u = 4.5$ . Let's see if mean residual life plot and parameter estimates against threshold graph confirm this choice.



**Figure 2.50:** Mean residual life plot (on the left) and parameter estimates against threshold (on the right) for "Islands" region.

Mean residual life plot is essentially linear until  $u = 5.2$ , but above this value there are too few observations to perform a model. Instead the other graph shows a big increase in sampling error over  $u = 4.5$ , thus this value can be chosen as threshold.

Shape parameter estimate for  $u = 4.5$  isn't significant as we can see from the following table.

```

ML fit of generalized Pareto distribution          Number of obs =          90
Log likelihood = -34.513473                      Wald chi2(0) =           .
                                                Prob > chi2 =            .

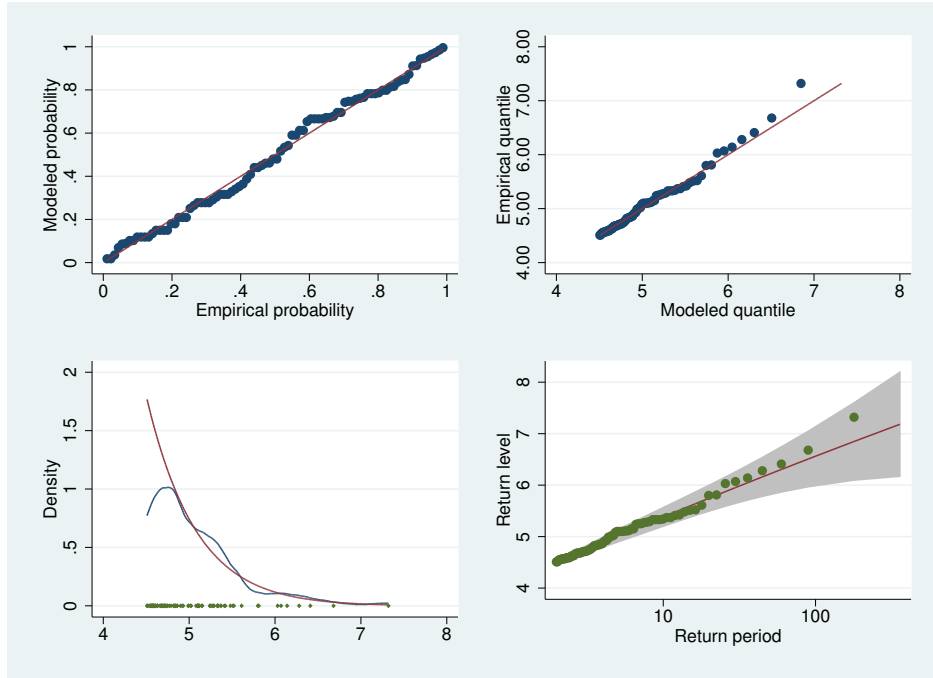
```

MwDef	Coef.	Std. Err.	z	P> z	[95% Conf. Interval]
lnsig					
_cons	-.5870707	.1508826	-3.89	0.000	-.8827952    -.2913463
xi					
_cons	<b>-.0294462</b>	.1080004	-0.27	<b>0.785</b>	<b>-.2411231    .1822306</b>
sig	<b>.5559534</b>	.0838837			.4136251    .7472569

**Figure 2.51:** Parameter estimates and 95% confidence intervals for  $u = 4.5$ .

As the "South" region case, if we try to fit a model with  $\xi = 0$  the p-value obtained for scale parameter estimate is one and AIC and BIC values for the unrestricted model are much lower, indicating a better fit.

However model diagnostic plots suggest model adequacy.



**Figure 2.52:** Probability plot (on the top left) and Quantile plot (on the top right), Kernel density plot (on the bottom left) and Return level plot (on the bottom right) for  $u = 4.5$ .

### 2.4.3 Results comparison

With the aim to summarize results obtained in previous subsections, we recall them in a table<sup>3</sup>.

Zone	$u$	$\hat{\xi}$	95% C.I.	Dist.Fam.
Italy	4.5	-0.1068006 **	[-0.1745711 , -0.03903]	Weibull
"Sea/Foreign" region	4.6	-0.3496852 *	[-.6269934 , -.072377]	Weibull
"Alps" region	4.4	-0.1804924 *	[-0.3203359 , -0.0406489]	Weibull
"Po valley" region	4.6	-0.2025167 *	[-0.3581311 , -0.0469023]	Weibull
"Centre" region	4.5	-0.2135829 ***	[-0.3239422 , -0.1032236]	Weibull
"South" region	4.6	-0.1068119	[-0.3716103 , 0.1579866]	Weibull
"Islands" region	4.5	-0.0294462	[-0.2411231 , 0.1822306]	Weibull

<sup>3</sup>one star (\*) if  $p < 0.05$ , two stars (\*\*) if  $p < 0.01$ , and three stars (\*\*\*) if  $p < 0.001$ .

### 3.1 Introduction

Natural disasters as earthquakes pose several challenges to insurers because they involve potentially high losses that are extremely uncertain. Before insurance providers are willing to offer coverage against an uncertain event two conditions are to be satisfied:

1. being able to identify and quantify, or estimate partially, the probabilities that an event occurs and the extent of losses likely to be incurred;
2. being able to set premiums for each potential customer.

The first condition can be satisfied estimating probabilities of specific events and the likely extent of losses using past data and catastrophe models (i.e. those available in Risk Management Solutions software); the second one can be satisfied charging premiums until better estimates of the risk are available.

### 3.2 Recent catastrophic earthquakes

Every year Italian seismographs detect thousands earthquakes, fortunately only few of them are catastrophic.

From 1900 to 2014 in the catalogue CPTI15 there are 36 earthquakes with magnitude greater or equal to 5.64, ten of which so disastrous as to be included in the list of "emergencies" by the Italian Civil Protection.

Year	Place	MwDef	Victims	Evacuees
1908	Messina	7.1	86000	missing value
1915	Avezzano	7.08	30000	missing value
1930	Vulture	6.67	1404	missing value
1968	Belice	6.41	296	90000
1976	Friuli	6.45	965	45000
1980	Irpinia	6.81	2734	70000
1997	Umbria-Marche	5.97	11	80000
2002	Molise	5.92	30	3000
2009	Abruzzo	6.29	308	80000
2012	Emilia	6.09	27	15000
2016	Amatrice	6.00	295 (at least)	4700 (at least)

*Remark 3.2.1.* The red row refers to the recent earthquake with epicentre in Accumuli (Rieti) and data are not included in the catalogue CPTI15 but were taken from the INGV site (see ING [19]).

Italian high seismic leads experts developing methods for evaluating damages and their costs: Civil Protection releases to highly qualified person forms which allow to carry out a damage survey on the whole building heritage hit by the seism in an homogeneous manner. Often, however, costs estimates refer only to structural damages and not to those related to the interruption of activities; so the Italian system is not completely adequate.

For 7 of the 10 events listed above, being concentrated in a little period of time (1968-2014), is possible obtaining data economically comparable. See CNI [21].

Year	Place	MwDef	Cost (billion Euro)
1968	Belice	6.41	9.179(*)
1976	Friuli	6.45	18.54(*)
1980	Irpinia	6.81	52.026(**)
1997	Umbria-Marche	5.97	13.463(*)
2002	Molise	5.92	1.4(*)
2009	Abruzzo	6.29	13.7
2012	Emilia	6.09	13.3

Discounted cost (\*) 2014 (\*\*) 2008. Costs are based on public funding except for Abruzzo and Emilia, which costs are expenditure forecasts.

### 3.3 Insurer's loss due to catastrophic events

Let us assume that an insurer insures for seismic risk the totality of Italian residential buildings and in particular he insures only the maximum event recorded in a year.

As we know from the previous chapter, the Moment Magnitude scale is defined by

$$M_w = \frac{2}{3}(\log_{10}M_0 - 6.03)$$

where  $M_0$  is the seismic moment at the focus measured in  $N \cdot m$ , in practice an energy; thus we will use the symbol  $E$  to denote it instead of  $M_0$ .

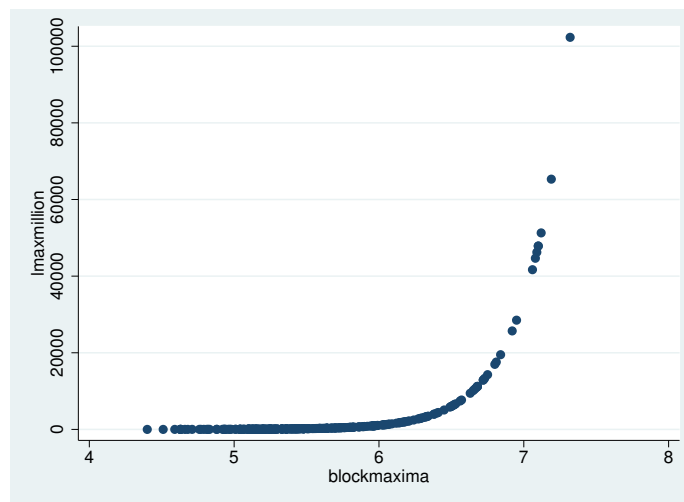
The energy released from an earthquake is strictly related with its destroying power, so we can assume that energy is directly proportional to the damage and in particular to the extent of loss, that is

$$L = kE = k10^{\frac{3}{2}M_w+6.03}. \quad (3.1)$$

The magnitude ( $m_w$ ) and cost ( $l$ ) values of the 7 recent catastrophic earthquakes can be used to find an approximate value for the proportionality constant  $k$ :

$$k = \frac{1}{7} \sum_{i=1}^7 \frac{l_i}{10^{\frac{3}{2}m_{w_i}+6.03}} = 5.449 \cdot 10^{-6} \approx 10^{-6}.$$

Using equation (3.1), we infer the extent of loss for all annual block maxima.



**Figure 3.1:** Loss (million Euro) against magnitude for all block maxima.

So in our context  $L$  is a random variable depending on value of annual block maxima  $M_w$ .

In the second chapter (Section 2.3.1) we found that the distribution of Italian annual block maxima is well represented by the Gumbel distribution, which estimated parameters are the following.

```

ML fit of generalized extreme value distribution   Number of obs   =       270
                                                    Wald chi2(0)    =       .
Log likelihood = -233.97668                       Prob > chi2     =       .
( 1)  [xi]_cons = 0

```

blockmaxima	Coef.	Std. Err.	z	P> z	[95% Conf. Interval]	
mu						
_cons	5.362347	.0318642	168.29	0.000	5.299895	5.4248
lnsig						
_cons	-.7003485	.0471259	-14.86	0.000	-.7927135	-.6079834
xi						
_cons	0 (omitted)					
sig	.4964123	.0233939			.452615	.5444477

**Figure 3.2:** Parameter estimates for annual block maxima with constraint  $\xi = 0$ .

Substituting estimates in the GEV distribution function formula

$$H_{0;\mu,\sigma}(x) = \exp \left\{ -e^{-\left(\frac{x-\mu}{\sigma}\right)} \right\}$$

and subtracting the found value from 1, we obtain the exceedance probability for a specified value of magnitude. Its inverse is the relative return period.

The expected annual maximum value of magnitude for the estimated Gumbel distribution is

$$\mathbb{E}[X] = \hat{\mu} + \gamma\hat{\sigma} = 5.362347 + 0.57722 \cdot 0.4964123 = 5.6488861$$

and its standard deviation is

$$StDev = \sqrt{Var(X)} = \frac{\pi\hat{\sigma}}{\sqrt{6}} = \pi \frac{0.4964123}{\sqrt{6}} = 0.63635074$$

leading to the upper bound value  $\mathbb{E}[X] + 1 \cdot StDev = 6.2852367$ .  
For definitions and formulas see Section 1.4.

For the expected annual maximum value of magnitude found above, that is 5.6488861, we have loss equal to 318.66117 million Euro, exceedance probability equal to 0.42962268 and return period of 2.3276239 years, which has to be interpreted as the minimum amount of time necessary for exceeding the magnitude value 5.6488861.

These values are coherent with those calculated by the software for a magnitude value very close to the expected annual maximum:

blockm~a	probde~y	lmaxmi~n	exprob	retper
5.65	.6444785	319.8896	.4289046	2.331521

**Figure 3.3:** Probability density, loss (million Euro), exceedance probability and return period (years) for magnitude 5.65.

The median can be found substituting  $p = 0.5$  in the quantile function for Gumbel distribution (equation (1.4))

$$\hat{x}_{0.5} = \hat{\mu} - \hat{\sigma} \ln(-\ln(0.5)) = 5.5442885$$

with a probability density value 0.69815674 and loss 222.04078 million Euro.

The modal value instead could be obtained equating to zero the estimated Gumbel density function derivative

$$h'_{0;\hat{\mu},\hat{\sigma}}(x) = \frac{H_{0;\hat{\mu},\hat{\sigma}}(x)e^{-\frac{x-\hat{\mu}}{\hat{\sigma}}}[-1 + e^{-\frac{x-\hat{\mu}}{\hat{\sigma}}}]}{\hat{\sigma}^2},$$

obtaining

$$\hat{x}_{modal} = \hat{\mu} = 5.362347$$

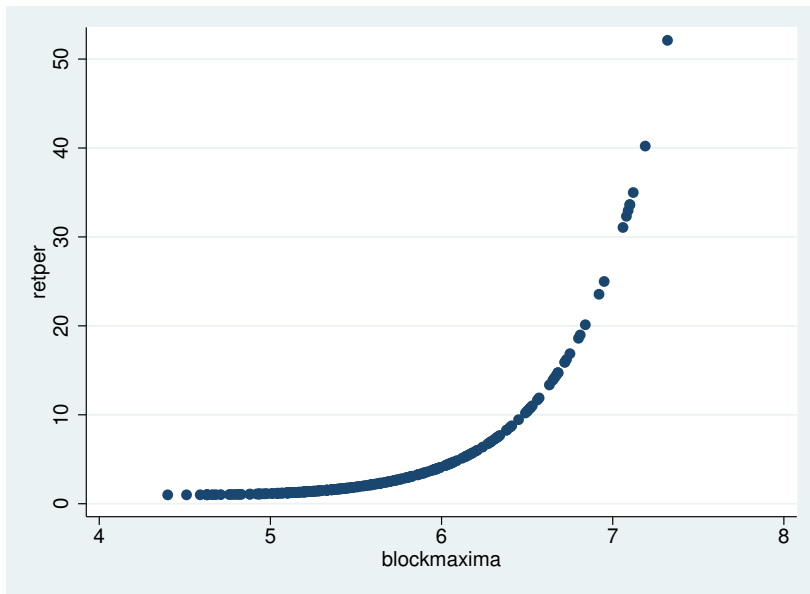
associated to a probability density value 0.7410764, a loss of 117.4898 million Euro, an exceedance probability of 0.63212056 and a return period of 1.5819767 years, all coherent with values computed by the software

blockm~a	probde~y	lmaxmi~n	exprob	retper
5.36	.7410681	117.4898	.6338601	1.577635

**Figure 3.4:** Probability density, loss (million Euro), exceedance probability and return period (years) for magnitude 5.36.

Obviously, return period increases increasing the magnitude:





**Figure 3.5:** Estimated return periods against magnitude for all block maxima.

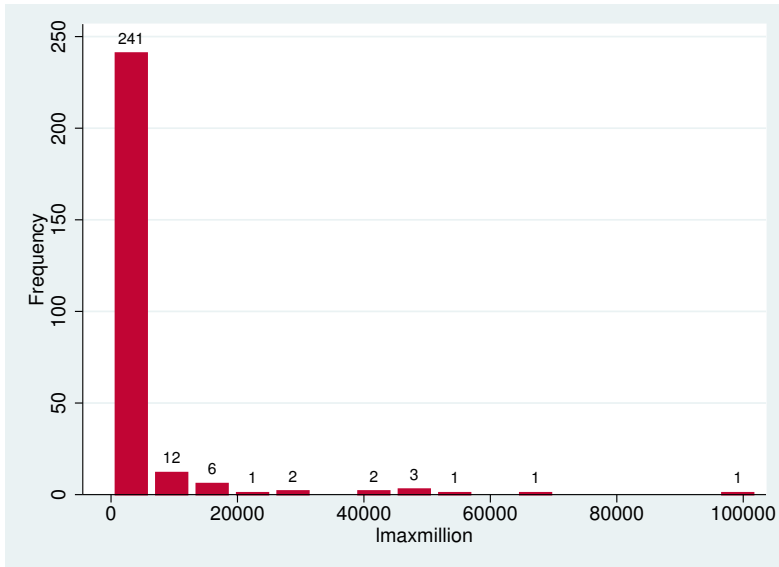
Place	MwDef ( <i>magnitude</i> )	Exceedance probability $1 - p$	Return period $1/(1 - p)$ (years)
Messina	7.1	0.0297336	33.63195
Avezzano	7.08	0.0309369	32.32383
Irpinia	6.81	0.052698	18.97603
Vulture	6.67	0.0692605	14.43824
Friuli	6.45	0.1057786	9.453704
Belice	6.41	0.1141287	8.762035
Abruzzo	6.29	0.143004	6.992813
Emilia	6.09	0.2061715	4.850332
Umbria-Marche	5.97	0.2547427	3.92553
Molise	5.92	0.2776051	3.602239

Is important underlying that the distribution of  $L$  isn't the same of the annual maximum magnitude, but its transformation. In fact, having the annual block maxima  $M_w$  distribution function  $H_{0;5.362347,0.4964123}$ ,  $L$  has distribution function

$$\mathbb{P}(L \leq \lambda) = \mathbb{P}(k10^{\frac{3}{2}M_w+6.03} \leq \lambda)$$

$$\begin{aligned}
&= \mathbb{P}\left(10^{\frac{3}{2}M_w} \leq \frac{\lambda}{k10^{6.03}}\right) \\
&= \mathbb{P}\left(\frac{3}{2}M_w \leq \log_{10} \frac{\lambda}{k10^{6.03}}\right) \\
&= \mathbb{P}\left(M_w \leq \frac{2}{3} \log_{10} \frac{\lambda}{k10^{6.03}}\right) \\
&= H_{0;5.362347,0.4964123}\left(\frac{2}{3} \log_{10} \frac{\lambda}{k10^{6.03}}\right).
\end{aligned}$$

Looking to frequencies we can think that  $L$ 's distribution is well represented by the Fréchet distribution.



**Figure 3.6:** Frequencies of loss  $L$  (million Euro).

Trying to fit a GEV model we obtain the table below, where we can see that estimates are all highly significant and  $\hat{\xi} = 1.678477 > 0$ , confirming that Fréchet distribution represents data well.

Remember that for  $\xi \geq 1$  the expected value for the GEV distribution is infinity, and for  $\xi \geq 1/2$  also its variance.

Anyway in practice the loss has an upward limit: in fact reconstruction costs of the whole Italian building heritage, composed by 27 million housing units, amount of 3900 billion Euro (see ANI [22] ).

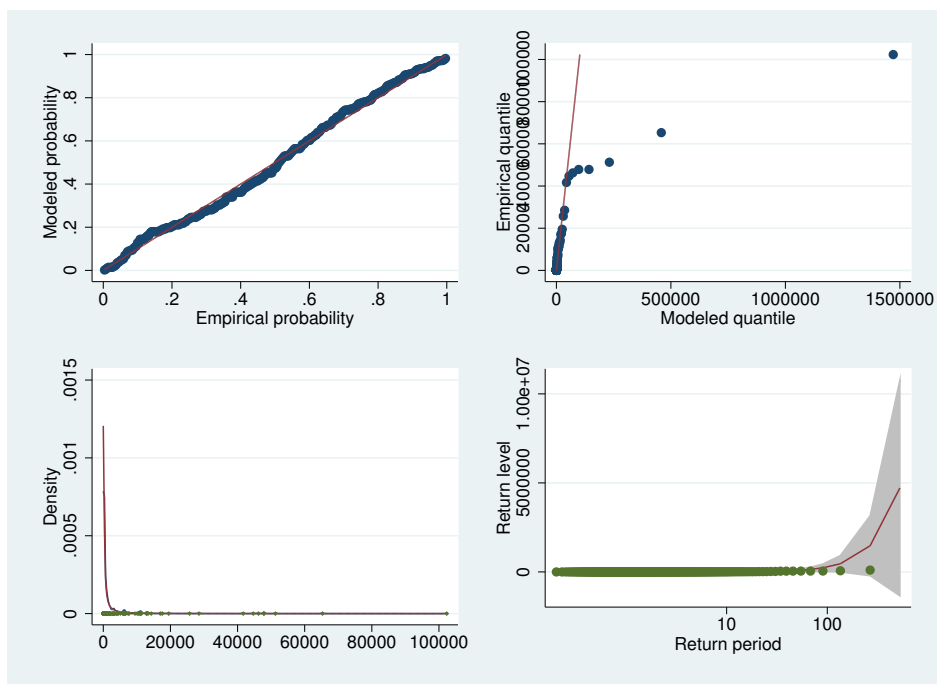
```

ML fit of generalized extreme value distribution   Number of obs   =       270
                                                Wald_chi2(0)      =       .
Log likelihood = -2120.0434                    Prob > chi2       =       .

```

lmaxmillion		Coef.	Std. Err.	z	P> z	[95% Conf. Interval]
mu	_cons	120.508	13.67867	8.81	0.000	93.69828 147.3177
	lnsig	5.319494	.1313234	40.51	0.000	5.062105 5.576883
xi	_cons	1.678477	.1035159	16.21	0.000	1.47559 1.881365
	sig	204.2805	26.82682			157.9226 264.2468

**Figure 3.7:** Parameter estimates for loss values (million Euro).



**Figure 3.8:** Probability plot and Quantile plot (on the top), Kernel density plot and Return level plot (on the bottom) for the fitted model.

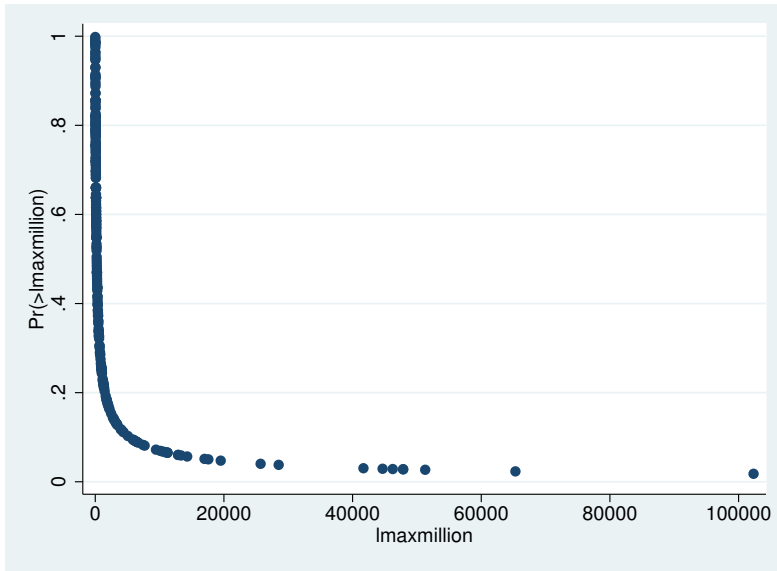
Diagnostic plots seems to confirm model adequacy.

Since the loss distribution function is the Fréchet, given by

$$H_{\xi;\mu,\sigma}(\lambda) = \exp \left\{ - \left[ 1 + \xi \left( \frac{\lambda - \mu}{\sigma} \right) \right]^{-1/\xi} \right\}$$

where  $\xi = 1.678477$ ,  $\mu = 120.508$  and  $\sigma = 204.2805$ , we can now calculate exceedance probability and return period for each value of loss.

After that we can plot exceedance probabilities against their loss values obtaining a graph called **exceedance probability curve**.



**Figure 3.9:** Exceedance probability curve.

The insurer can use this graph to determine how large a loss will occur at a given probability level: the value obtained is called **probable maximum loss**. Often the PML is associated to its return period instead of its exceedance probability.

The exceedance probability curve can also be used to distribute loss: an example is a homeowner having deductible on his insurance policy such that he had to cover the first part of the loss, an insurer covers the middle portion and a reinsurer handles the loss above a certain amount.

### 3.4 Premium and equilibrium reserve

In the past section we saw that annual maximum magnitude  $M_w$  and loss  $L$  distribution are related thus, for simplicity, from now on we will refer to  $M_w$ 's distribution.

To begin can be useful give a look to annual maximum magnitude quantiles, given by

$$\hat{x}_p = \hat{\mu} - \hat{\sigma} \ln(-\ln p) \quad \text{if } \hat{\xi} = 0$$

where  $\hat{\mu} = 5.362347$  and  $\hat{\sigma} = 0.4964123$ , and their return periods (in years). Thus we have

$$\begin{aligned} \hat{x}_{0.01} &= \hat{\mu} - \hat{\sigma} \ln(-\ln 0.01) = 4.6042362 & \text{return period} &= 1.010101 \\ \hat{x}_{0.05} &= \hat{\mu} - \hat{\sigma} \ln(-\ln 0.05) = 4.817689 & \text{return period} &= 1.0526316 \\ \hat{x}_{0.10} &= \hat{\mu} - \hat{\sigma} \ln(-\ln 0.10) = 4.948323 & \text{return period} &= 1.1111111 \\ \hat{x}_{0.25} &= \hat{\mu} - \hat{\sigma} \ln(-\ln 0.25) = 5.2002017 & \text{return period} &= 1.3333333 \\ \hat{x}_{0.50} &= \hat{\mu} - \hat{\sigma} \ln(-\ln 0.50) = 5.5442885 & \text{return period} &= 2 \\ \hat{x}_{0.75} &= \hat{\mu} - \hat{\sigma} \ln(-\ln 0.75) = 5.9808267 & \text{return period} &= 4 \\ \hat{x}_{0.90} &= \hat{\mu} - \hat{\sigma} \ln(-\ln 0.90) = 6.479457 & \text{return period} &= 10 \\ \hat{x}_{0.95} &= \hat{\mu} - \hat{\sigma} \ln(-\ln 0.95) = 6.8367885 & \text{return period} &= 20 \\ \hat{x}_{0.99} &= \hat{\mu} - \hat{\sigma} \ln(-\ln 0.99) = 7.6459177 & \text{return period} &= 100 \end{aligned}$$

which are pretty close to empirical quantiles, confirming another time appropriateness of the estimated Gumbel distribution.

Also the empirical mean is very similar to the expected annual maximum value estimated in Section 3.3.

blockmaxima				
	Percentiles	Smallest		
1%	4.59	4.4		
5%	4.81	4.51		
10%	5.01	4.59	Obs	270
25%	5.19	4.63	Sum of Wgt.	270
50%	5.51		Mean	5.643778
		Largest	Std. Dev.	.6010701
75%	6.03	7.1		
90%	6.525	7.12	Variance	.3612853
95%	6.75	7.19	Skewness	.5972236
99%	7.12	7.32	Kurtosis	2.766683

Figure 3.10: Empirical quantiles for annual maximum magnitude value.

Furthermore, to each quantile value we can associate its loss

<b>MwDef</b> ( <i>magnitude</i> )	<b>Energy</b> ( $N \cdot m$ )	<b>Loss</b> ( <i>million Euro</i> )	<b>Exceedance Prob.</b> ( $1 - p$ )	<b>Return Period</b> ( <i>years</i> )
4.6042362	8.637e+12	8.6368286	0.99	1.010101
4.817689	1.805e+13	18.05234	0.95	1.0526316
4.948323	2.835e+13	28.345525	0.90	1.1111111
5.2002017	6.766e+13	67.655413	0.75	1.3333333
5.5442885	2.220e+14	222.04078	0.50	2
5.9808267	1.003e+15	1002.8594	0.25	4
6.479457	5.613e+15	5612.8767	0.10	10
6.8367885	1.928e+16	19283.362	0.05	20
7.6459177	3.154e+17	315410.79	0.01	100

an then calculate the expected annual loss

$$EAL = \sum_{i=1}^9 (1 - p_i) L_i = 5143.2516 \text{ million Euro.}$$

Keeping in mind that in the whole Italian territory there are 27 million housing units with an overall reconstruction cost of 3900 billion Euro, we can calculate the premium assuming an obligatory and solidarity insurance system:

$$\frac{EAL}{\text{num. of housing units}} = \frac{5143.2516 \cdot 10^6}{27 \cdot 10^6} = 190.49 \text{ Euro}$$

for a house of  $3900 \cdot 10^9 / 27 \cdot 10^6 \approx 144444$  Euro.

Thus, for an housing unit with a reconstruction value of 100000 Euro we obtain

$$\text{premium} \approx 131.88 \text{ Euro.}$$

Premiums collected every year by the insurer are needed for covering the expected annual loss; obviously the real annual loss is a priori unknown and its extent could be greater or less than the expected.

This implies that, in years in which the real loss is less than the expected one, the insurer sets aside part of collected premiums earmarking them to the coverage of future losses greater than the expected one. This reserve is called

**equilibrium reserve:** it can be established when the **loss ratio**, given by the ratio between real loss and collected premiums, is much less than 1 and can be used when the ratio is greater than 1.

Being a reserve and not a profit for the insurer, this amount of money should be not taxed.

*Remark 3.4.1.* Countries which adopted this measure first are Germany (1952) and Finland (1953); in Italy the equilibrium reserve for natural disasters should be established for all non-life branches, excluding the Credit and Bail branch which has a different kind of reserve (see Donati and Putzolu [23]).

### 3.5 Premium diversification among zones

Let us suppose now that housing units distribution over the Italian territory is homogeneous and we take in consideration only populated zones, that is excluding the *Sea/Foreign* region.

Under this hypothesis there are 5 zones each with 5.4 million housing units with an overall reconstruction value of 780 billion Euro.

Using the quantile formula (1.4)

$$\hat{x}_p = \begin{cases} \hat{\mu} + [(-\ln p)^{-\hat{\xi}} - 1]\hat{\sigma}/\hat{\xi} & \text{if } \hat{\xi} \neq 0 \\ \hat{\mu} - \hat{\sigma} \ln(-\ln p) & \text{if } \hat{\xi} = 0 \end{cases}$$

we calculate quantiles and losses for each zone and then their respective premiums<sup>1</sup>.

- *Alps* region:  $\hat{\xi} = 0$ ,  $\hat{\mu} = 4.592146$ ,  $\hat{\sigma} = 0.4222104$

MwDef ( <i>magnitude</i> )	Loss ( <i>million Euro</i> )	Exceedance Prob. ( $1 - p$ )	Return Period ( <i>years</i> )
3.9473549	0.89337203	0.99	1.010101
4.1289015	1.6724466	0.95	1.0526316
4.2400088	2.4547835	0.90	1.1111111
4.4542376	5.1446567	0.75	1.3333333
4.7468916	14.136354	0.50	2
5.1181777	50.964357	0.25	4
5.5422745	220.5016	0.10	10
5.8461933	629.9266	0.05	20
6.5343768	6785.2398	0.01	100

<sup>1</sup>Parameter values used are obtained fitting a GEV model with the constraint  $\xi = 0$ , except for the *Centre* region whose estimates are those found in Subsection 2.3.4.

The expected annual loss is

$$EAL = \sum_{i=1}^9 (1 - p_i) L_i = 149.74921 \text{ million Euro.}$$

Dividing this value for the 5.4 million housing units in the region and relating it to an house value of 100000 Euro we obtain

$$\text{premium} \approx 19.20 \text{ Euro}$$

- *Po valley* region:  $\hat{\xi} = 0$ ,  $\hat{\mu} = 4.760912$ ,  $\hat{\sigma} = 0.4672812$

<b>MwDef</b> ( <i>magnitude</i> )	<b>Loss</b> ( <i>million Euro</i> )	<b>Exceedance Prob.</b> ( $1 - p$ )	<b>Return Period</b> ( <i>years</i> )
4.0472897	1.2616374	0.99	1.010101
4.2482163	2.5253667	0.95	1.0526316
4.3711843	3.8616815	0.90	1.1111111
4.608282	8.7583642	0.75	1.3333333
4.9321766	26.80803	0.50	2
5.3430973	110.8271	0.25	4
5.8124663	560.65984	0.10	10
6.1488284	1791.6064	0.05	20
6.9104753	24872.128	0.01	100

The expected annual loss is

$$EAL = \sum_{i=1}^9 (1 - p_i) L_i = 449.17078 \text{ million Euro.}$$

Dividing this value for the 5.4 million housing units in the region and relating it to an house value of 100000 Euro we obtain

$$\text{premium} \approx 57.59 \text{ Euro}$$

- *Centre* region:  $\hat{\xi} = -0.1371585$ ,  $\hat{\mu} = 4.966685$ ,  $\hat{\sigma} = 0.5771983$



<b>MwDef</b> ( <i>magnitude</i> )	<b>Loss</b> ( <i>million Euro</i> )	<b>Exeedance Prob.</b> ( $1 - p$ )	<b>Return Period</b> ( <i>years</i> )
3.9860809	1.0212248	0.99	1.010101
4.2832546	2.8502431	0.95	1.0526316
4.4566673	5.1880117	0.90	1.1111111
4.7738653	15.516649	0.75	1.3333333
5.1730063	61.589893	0.50	2
5.6277268	296.20351	0.25	4
6.0842609	1433.479	0.10	10
6.3748257	3910.5564	0.05	20
6.9357863	27144.35	0.01	100

The expected annual loss is

$$EAL = \sum_{i=1}^9 (1 - p_i) L_i = 735.19049 \text{ million Euro.}$$

Dividing this value for the 5.4 million housing units in the region and relating it to an house value of 100000 Euro we obtain

$$premium \approx 94.26 \text{ Euro}$$

- *South* region:  $\hat{\xi} = 0$ ,  $\hat{\mu} = 4.889862$ ,  $\hat{\sigma} = 0.6290228$

<b>MwDef</b> ( <i>magnitude</i> )	<b>Loss</b> ( <i>million Euro</i> )	<b>Exeedance Prob.</b> ( $1 - p$ )	<b>Return Period</b> ( <i>years</i> )
3.9292312	0.83916391	0.99	1.010101
4.1997053	2.1357871	0.95	1.0526316
4.3652366	3.7831617	0.90	1.1111111
4.6844016	11.392064	0.75	1.3333333
5.120407	51.358283	0.50	2
5.6735611	347.00973	0.25	4
6.3053944	3076.7457	0.10	10
6.7581825	14698.525	0.05	20
7.7834607	507213.79	0.01	100

The expected annual loss is

$$EAL = \sum_{i=1}^9 (1 - p_i) L_i = 6241.9789 \text{ million Euro.}$$

Dividing this value for the 5.4 million housing units in the region and relating it to an house value of 100000 Euro we obtain

$$premium \approx 800.26 \text{ Euro}$$

- *Islands* region:  $\hat{\xi} = 0$ ,  $\hat{\mu} = 4.66749$ ,  $\hat{\sigma} = 0.4133187$

<b>MwDef</b> ( <i>magnitude</i> )	<b>Loss</b> ( <i>million Euro</i> )	<b>Exceedance Prob.</b> ( $1 - p$ )	<b>Return Period</b> ( <i>years</i> )
4.0362781	1.2145549	0.99	1.010101
4.2140014	2.2438928	0.95	1.0526316
4.3227688	3.2670306	0.90	1.1111111
4.532486	6.7410725	0.75	1.3333333
4.8189766	18.132801	0.50	2
5.1824435	63.630487	0.25	4
5.5976089	266.93981	0.10	10
5.8951272	745.91751	0.05	20
6.5688177	7642.3434	0.01	100

The expected annual loss is

$$EAL = \sum_{i=1}^9 (1 - p_i) L_i = 176.71755 \text{ million Euro.}$$

Dividing this value for the 5.4 million housing units in the region and relating it to an house value of 100000 Euro we obtain

$$premium \approx 22.66 \text{ Euro.}$$

Also in this case, since the extent of the real annual loss could be greater or less than the expected, we can establish an equilibrium reserve with the purpose of compensating the fluctuations in time of the loss.

## CONCLUSIONS

In this thesis we addressed the extreme value theory focusing on the distribution of maxima.

The most important result is the Fisher-Tippet theorem, which identifies three possible distribution functions families for the normalised maxima, that is the Gumbel, the Fréchet and the Weibull distribution function.

The theory described in the first chapter was then applied to Italian earthquakes data available in the Parametric Catalogue of Italian Earthquakes (CPTI15), taking into account maxima values of the recorded moment magnitude.

From the analysis we obtained maximum likelihood estimates of the distribution function parameters, first for the whole Italian territory and then for each seismic macro-zone; the adequacy of each fitted model was then checked using probability plots, quantile plots, return level plots and kernel density plots.

Information on the most suitable distribution function for the maxima let us doing some actuarial evaluations: we calculated the insurer's expected annual loss and then the premium that people have to pay if they want to insure homes from seismic risk.

In particular we saw that the extent of premium differs region by region depending on zone seismic. This difference open the way for two possibilities: set a single premium for all zones or differentiate it depending on the taken risk.

Actually in Italy there isn't the obligation of insurance coverage against natural disasters, even if there's the possibility of extending fire policies to natural disaster damages. This fact leads to an high anti-selected request that does not allow insurers to offer the same price on the whole Italian territory: in fact citizens more sensitive to the insurance coverage are probably those who

live in most risky zones.

In this case the only way could be making compulsory the coverage, raising however the problem of the high capital requirements needed to ensure the solvency of the insurance company. In such situation becomes necessary introducing overdrawn and/or deductibles or a public reinsurer.

Instead, if the insurer provides a zone differentiated premium probably many more citizens insure their homes, but those who need it most should pay a too high premium: in this case could be useful introducing tax incentives and exempting from tax payment (in Italy equal to 22.25%, one of the highest in Europe), especially if the coverage is not compulsory.

## BIBLIOGRAPHY

- [1] Paul Embrechts, Claudia Kluppelberg, and Thomas Mikosch. *Modelling Extremal Events for Insurance and Finance*. Springer, 1997.
- [2] Sidney I. Resnick. *Extreme Values, Regular Variation and Point Processes*. Springer-Verlag, 1987.
- [3] J. R. M. Hosking and James R. Wallis. *Regional Frequency Analysis: An Approach Based on L-moments*. Cambridge University Press, 1997.
- [4] Samuel Kotz and Saralees Nadarajah. *Extreme Value Distributions: Theory and Applications*. Imperial College Press, 2000.
- [5] Rolf-Dieter Reiss and Michael Thomas. *Statistical Analysis of Extreme Values: from Insurance, Finance, Hydrology and Other Fields*. Springer Basel AG, 1997. ISBN 978-3-7643-5768-9.
- [6] Stuart Coles. *An Introduction to Statistical Modeling of Extreme Values*. Springer, 2001.
- [7] Anthony C. Davison. *Statistical Models*. Cambridge University Press, 2003.
- [8] Asuka Suzuki-Parker. *Uncertainties and Limitations in Simulating Tropical Cyclones*. Springer-Verlag Berlin Heidelberg, 2012.
- [9] Jan Beirlant, Yuri Goegebeur, Johan Segers, Jozef Teugels, Daniel De Waal, and Chris Ferro. *Statistics of Extremes: Theory and Applications*. John Wiley & Sons, 2004.

- [10] Emil J. Gumbel. *Statistics of Extremes*. Columbia University Press, 1958.
- [11] Ross Leadbetter, Stamatis Cambanis, and Vladas Pipiras. *A Basic Course in Measure and Probability: Theory for Applications*. Cambridge University Press, 2014.
- [12] Sheri M. Markose and Amadeo Alentorn. The generalized extreme value (gev) distribution, implied tail index and option pricing. *University of Essex Department of Economics Discussion Papers*, 2005.
- [13] Hiroo Kanamori. Earthquake seismology. In Gerald Schubert, editor, *Treatise on Geophysics*. Elsevier, 2007.
- [14] Istituto Nazionale di Geofisica e Vulcanologia - Sezione di Bologna, . URL <http://www.bo.ingv.it/italiano/ricerca/Sismologia/Macrosismica-e-Sismologia-Storica/>.
- [15] Luigi Vannucci. Assicurazione e riassicurazione delle catastrofi ambientali. In Sara Landini and Giampiero Maracchi, editors, *Cambiamenti climatici, catastrofi ambientali e assicurazione*, pages 27–43. Fondazione CESIFIN Alberto Predieri, 2016. ISBN 978-88-98742-02-8.
- [16] Centro Geofisico Prealpino - Provincia di Varese. URL [http://www.astrogeo.va.it/sismologia/scale\\_sismiche.php](http://www.astrogeo.va.it/sismologia/scale_sismiche.php).
- [17] Andrea Rovida, Mario Locati, Romano Camassi, Barbara Lolli, and Paolo Gasperini (eds) 2016. Cpti15, the 2015 version of the parametric catalogue of italian earthquakes. *Istituto Nazionale di Geofisica e Vulcanologia*. doi:<http://doi.org/10.6092/INGV.IT-CPTI15>.
- [18] Istituto Nazionale di Geofisica e Vulcanologia - Sede Irpinia, . URL <http://www.gm.ingv.it/index.php/labgis/report-cartografici>.
- [19] Istituto Nazionale di Geofisica e Vulcanologia - Terremoti, . URL <https://ingvterremoti.wordpress.com/2016/08/24/evento-sismico-tra-le-province-di-rieti-e-ascoli-p-m-6-0-24-agosto/>.
- [20] Protezione Civile. URL [http://www.protezionecivile.gov.it/jcms/it/emerg\\_it\\_sismico.wp](http://www.protezionecivile.gov.it/jcms/it/emerg_it_sismico.wp).
- [21] Centro Studi Consiglio Nazionale Ingegneri - I costi dei terremoti in Italia (c.r 470). URL [https://www.tuttoingegnere.it/images/News/2016/I\\_costi\\_dei\\_terremoti\\_in\\_Italia.pdf](https://www.tuttoingegnere.it/images/News/2016/I_costi_dei_terremoti_in_Italia.pdf).

- [22] Associazione Nazionale fra le Imprese Assicuratrici - Danni da eventi sismici e alluvionali al patrimonio abitativo italiano: studio quantitativo e possibili schemi assicurativi. URL <http://www.ania.it/export/sites/default/it/pubblicazioni/monografie-e-interventi/Danni/Danni-da-eventi-sismici-e-alluvionali.pdf>.
- [23] Antigono Donati and Giovanna Volpe Putzolu. *Manuale di Diritto delle Assicurazioni*. Giuffrè Editore, 2012.
- [24] Luigi Vannucci. *Teoria del Rischio e Tecniche Attuariali Contro i Danni*. Pitagora Editrice, 2010.
- [25] Patricia Grossi, Howard Kunreuther, and Don Windeler. An introduction to catastrophe models and insurance. In Patricia Grossi and Howard Kunreuther, editors, *Catastrophe Modeling: A New Approach to Managing Risk*, pages 23–42. Springer, 2005. ISBN 0-387-23082-3.
- [26] Roberto Manzato. Aspetti quantitativi dei danni da catastrofi naturali e possibili opzioni di policy in materia assicurativa. In Sara Landini and Giampiero Maracchi, editors, *Cambiamenti climatici, catastrofi ambientali e assicurazione*, pages 27–43. Fondazione CESIFIN Alberto Predieri, 2016. ISBN 978-88-98742-02-8.

## *ACKNOWLEDGEMENTS*

First of all I would like to thank Prof. Luigi Vannucci for his support in writing this thesis and for the enthusiasm with which he embraced every my proposal.

Then I would like to thank the INGV (National Institute of Geophysics and Volcanology) and in particular Carlo Meletti for providing data necessary to my analysis and for his willingness.

Particular thanks go to Matteo for all good advice he gave me.

Finally I would like to thank my family for always believing in me.