

Extreme Value Theory & Time Series Analysis

Study year 2023/24, Q3/4

Fabian Mies

Delft University of Technology

f.mies@tudelft.nl

Structure of the course

- Part I: Extreme Value Theory in Q3
- Part II: Time Series Analysis in Q4
- each part has two mandatory assignments and an exam
- Lectures slides are put on Brightspace prior to the lecture

Structure of the course: assessment

Important dates:

- Extreme Value Theory, Assignment 1: hand in by March 05, 2024 (7.5% of overall grade)
- Extreme Value Theory, Assignment 2: hand in by March 26, 2024 (7.5% of overall grade)
- Extreme Value Theory, Exam: April 16, 2024 (30% of overall grade)
- Time Series Analysis, Assignment 1: hand in by May 14, 2024 (7.5% of overall grade)
- Time Series Analysis, Assignment 2: hand in by June 04, 2024 (7.5% of overall grade)
- Time Series Analysis, Exam: June 25, 2024 (30% of overall grade)

About the assignments

- Exercise sheet is published two weeks earlier
- Proofs and some programming / data analysis
- Hand-in via Brightspace as a group of 2

About the exam

- Written exam, 2h duration
- Contents: Proofs, routine calculations, interpretation
- You may bring 1 handwritten page as cheatsheet (A4, one side blank)

Extreme Value Theory & Time Series Analysis

Part I: Extreme Value Theory

Study year 2023/24, Q3

Fabian Mies

Delft University of Technology

f.mies@tudelft.nl

Outline

- | | |
|---|-------------|
| ① Introduction: Estimating the tail of a distribution | Lecture 1+2 |
| ② Generalized extreme value (GEV) distributions and the extreme value limit theorem | Lecture 2+3 |
| ③ Estimation of GEV parameters | Lecture 3+4 |
| ④ Assessing m-year returns | Lecture 4 |
| ⑤ Bivariate extremes | Lecture 5+6 |

Recommended literature:

de Haan & Ferreira (2006) Extreme value theory: an introduction ([ebook](#))

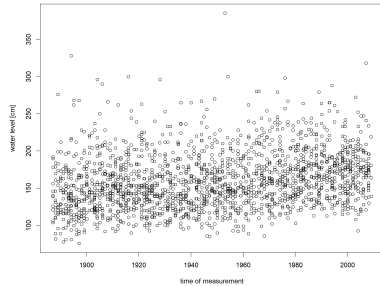
Outline

- | | | |
|---|---|-------------|
| ① | Introduction: Estimating the tail of a distribution | Lecture 1+2 |
| ② | Generalized extreme value (GEV) distributions and the extreme value limit theorem | Lecture 2+3 |
| ③ | Estimation of GEV parameters | Lecture 3+4 |
| ④ | Assessing m-year returns | Lecture 4 |
| ⑤ | Bivariate extremes | Lecture 5+6 |

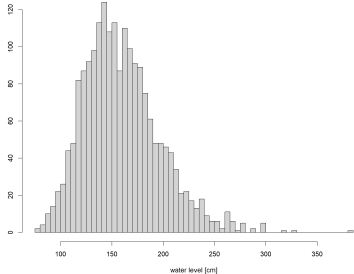
Recommended literature:

de Haan & Ferreira (2006) Extreme value theory: an introduction ([ebook](#))

Historical sea levels at Hoek van Holland



- ▶ Records of water levels during storms at Hoek van Holland, 1887-2009
- ▶ Measurements x_1, \dots, x_{1965} (multiple storms per year)

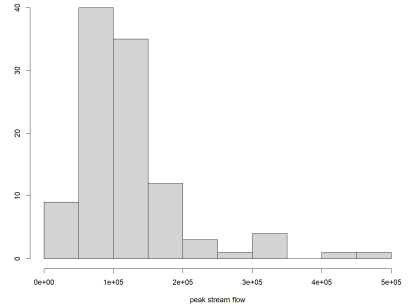
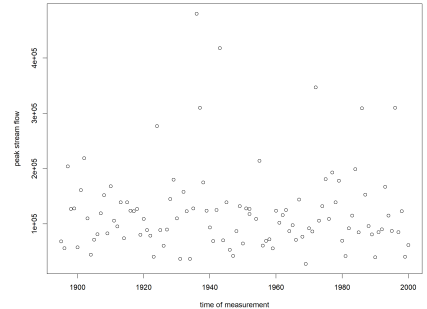


- ▶ How often should we expect a sea level above 400 cm, say?
- ▶ How high should the dike be to prevent 99.99% of all floods?

Peak stream flow of Potomac river



- ▶ Yearly maximum stream flow of Potomac river, from 1895-2000, measured at Point of Rocks, Maryland US
- ▶ Higher stream flow implies higher speeds and higher water levels



Estimating probabilities of rare events

For random variables $X_1, \dots, X_n \stackrel{iid}{\sim} F$ with cdf F , we may estimate $F(x) = P(X \leq x)$ by the empirical cdf (ecdf)

$$\hat{F}_n(x) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}(X_i \leq x).$$

Analogously, the exceedance probability $P(X > x) = 1 - F(x) = \bar{F}(x)$ can be estimated by the empirical survival function

$$\hat{\bar{F}}_n(x) = 1 - \hat{F}_n(x) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}(X_i > x).$$

Relative error of the ECDF

For any $x \in \mathbb{R}$ such that $F(x) \in (0, 1)$, we have

$$\frac{|\hat{\bar{F}}_n(x) - \bar{F}(x)|}{\bar{F}(x)} = \mathcal{O}_P \left(\frac{1}{\sqrt{n\bar{F}(x)}} \right).$$

Not reasonable for large x ,
i.e. small $\bar{F}(x)$.

For $\bar{F}(x) \ll \frac{1}{n}$, we
probably will not see ANY
data point exceeding x .

Example 1.1

The tail estimate can be better in parametric models:

- For some mean parameter $\mu \in \mathbb{R}$, we suppose that $X_1, \dots, X_n \stackrel{iid}{\sim} \mathcal{N}(\mu, 1)$.
- Estimate parameters by $\hat{\mu}_n = \frac{1}{n} \sum_{i=1}^n X_i$ such that

$$|\hat{\mu}_n - \mu| = \mathcal{O}_P(1/\sqrt{n}).$$

- Plug in the estimator to approximate the tail probability

$$P(X > x) = \bar{\Phi}(x - \mu) \approx \bar{\Phi}(x - \hat{\mu}_n)$$

- How does this approximation perform for large x ?

Relative error:

Lemma 1.2

For any $x > 1$, it holds $\frac{3}{4x} \leq \frac{\bar{\Phi}(x)}{\varphi(x)} \leq \frac{1}{x}$

Parametric vs Nonparametric approach

Nonparametric estimate:

- ▶ Estimate $P(X > x)$ by $\hat{\bar{F}}_n(x)$.
- ▶ Relative error for Gaussian observations ($\mu = 0$)

$$\begin{aligned}\frac{|\hat{\bar{F}}_n(x) - \bar{F}(x)|}{\bar{F}(x)} &= \mathcal{O}_P\left(\frac{1}{\sqrt{n\bar{\Phi}(x)}}\right) \\ &= \mathcal{O}\left(\frac{x}{\sqrt{n}} \exp\left(\frac{x^2}{2}\right)\right)\end{aligned}$$

Parametric estimate:

- ▶ Estimate $P(X > x)$ by $\bar{\Phi}(x - \hat{\mu}_n)$.
- ▶ Relative error for Gaussian observations ($\mu = 0$)

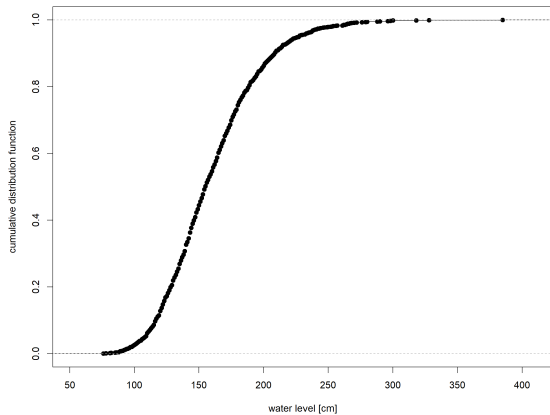
$$\frac{|\bar{\Phi}(x - \hat{\mu}_n) - \bar{F}(x)|}{\bar{F}(x)} = \mathcal{O}_P\left(\frac{x}{\sqrt{n}}\right)$$

Discuss

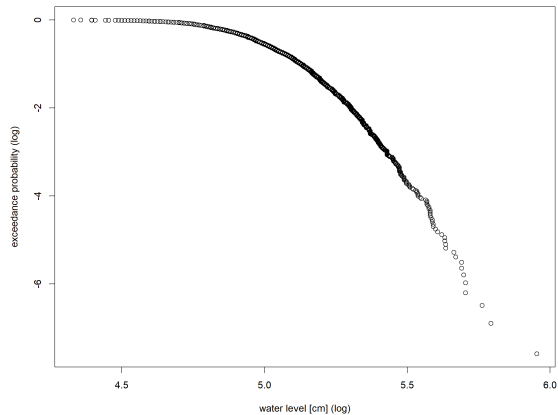
Which approach to use?

Hoek van Holland water levels (again)

Empirical CDF

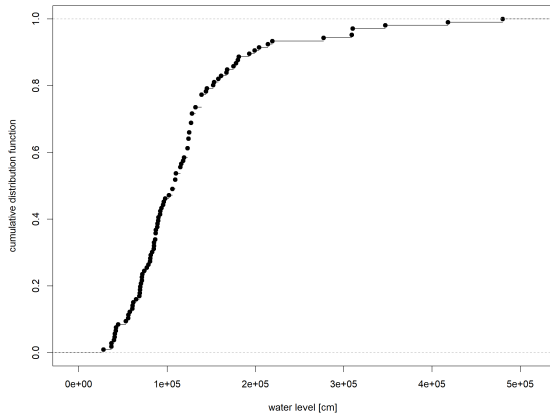


Log-log plot of empirical survival function

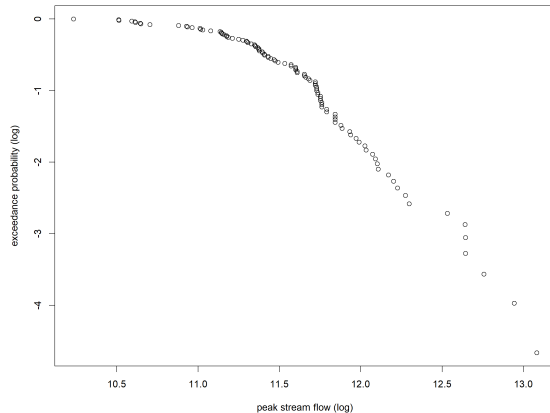


Potomac river stream flow (again)

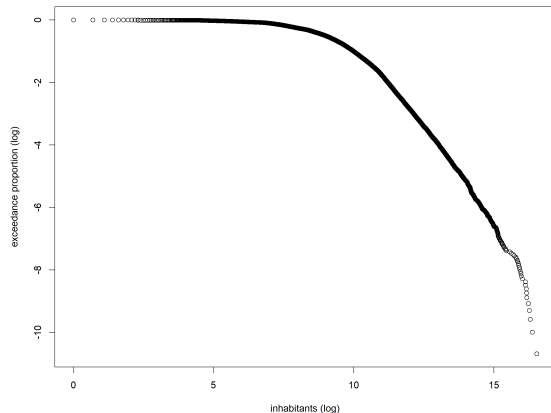
Empirical CDF



Log-log plot of empirical survival function



Log-log plot of empirical survival function



Findings

- log-log plot is approximately linear in the tail
- this suggests

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

for some $a > 0$ and $\alpha > 0$.

- The exponent $\gamma = 1/\alpha$ is referred to as the extreme value index of F
- **Central idea:** Use the data to find a and γ to extrapolate $\hat{\bar{F}}_n(x)$ for large values of x

Extreme value theory makes this extrapolation mathematically **rigorous**.

Making the log-log plot rigorous

Definition 1.3 (Regularly varying functions)

A function $f : (0, \infty) \rightarrow \mathbb{R}$ is regularly varying with index $\alpha \in \mathbb{R}$ if

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

If $\alpha = 0$, we say that f is slowly varying.

Exercise

Which functions are regularly varying?

- $f(x) = x^\alpha$
- $f(x) = \log(x)$
- $f(x) = x^\alpha \log(\log(x))$
- $f(x) = \exp(x)$

Lemma 1.4

Suppose that a function $f : (0, \infty) \rightarrow \mathbb{R}$ satisfies

$$\lim_{t \rightarrow \infty} \frac{f(tx)}{f(t)} = g(x), \quad \forall x > 0.$$

Then there exists some $\alpha \in \mathbb{R}$ such that $g(x) = x^\alpha$.

Proof:

Claim:

$$\text{If } \lim_{t \rightarrow \infty} \frac{f(tx)}{f(t)} = g(x)$$
$$\text{then } g(x) = x^\alpha$$

Proposition 1.5

Let X be a random variable with cdf F such that \bar{F} is regularly varying with exponent $-\frac{1}{\gamma}$ for some $\gamma > 0$. Then

$$P(X > tx | X > t) \xrightarrow{t \rightarrow \infty} \begin{cases} x^{-1/\gamma}, & x > 1, \\ 1, & x \leq 1. \end{cases}$$

That is, the conditional distribution of X/t given $X > t$ converges towards a Pareto distribution as $t \rightarrow \infty$.

We refer to γ as the **extreme value index** of F .

Definition 1.6 (Pareto distribution)

A random variable X has a Pareto distribution $X \sim \text{Par}(\alpha)$ with parameter $\alpha > 0$, if

$$P(X \leq x) = \max(0, 1 - x^{-\alpha}).$$

Different distributions with the same tail behavior

Definition 1.6 (Pareto distribution)

A random variable X has a Pareto distribution $X \sim \text{Par}(\alpha)$ with parameter $\alpha > 0$, if

$$P(X \leq x) = \max(0, 1 - x^{-\alpha}).$$

Definition 1.7 (Frechet distribution)

A random variable X has a Frechet distribution $X \sim \text{Frechet}(\alpha)$ with parameter $\alpha > 0$, if

$$P(X \leq x) = \exp(-x^{-\alpha}), \quad x > 0.$$

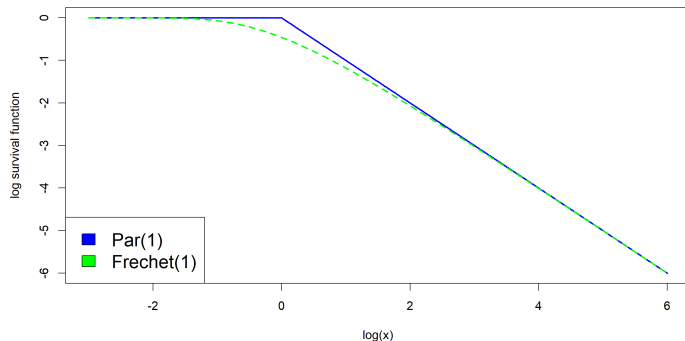


Figure 1: Theoretical log-log plot of the Pareto and Frechet distribution, with shape parameter $\gamma = 1$.

Exercise

Show that the Frechet distribution is regularly varying.

How to estimate the extreme value index γ of a distribution?

Example 1 (MLE for the Pareto distribution)

Let $X_1, \dots, X_n \stackrel{iid}{\sim} \text{Par}(1/\gamma)$ with unknown parameter γ .

- The Pareto pdf is

$$f_\gamma(x) = \frac{1}{\gamma} x^{-\frac{1}{\gamma}-1} \mathbb{1}(x > 1).$$

- Maximize the log-likelihood

$$l_n(X_1, \dots, X_n) = -n \log(\gamma) - \left(\frac{1}{\gamma} + 1\right) \sum_{i=1}^n \log(X_i)$$

$$\xrightarrow{\max} \hat{\gamma}_n = \frac{1}{n} \sum_{i=1}^n \log(X_i)$$

Exercise

Show that $\hat{\gamma}_n$ is consistent.

Semiparametric tail index estimation

Now let $X_1, \dots, X_n \stackrel{iid}{\sim} F$ for some cdf F with extreme value index $\gamma > 0$.

- ▶ Choose a threshold value $t > 0$ and only use data larger than t
- ▶ Treat those observations as if they were Pareto distributed with parameter $\frac{1}{\gamma}$:

$$\hat{\gamma}_n = \hat{\gamma}_n(t) = \frac{1}{n(t)} \sum_{i=1}^n \mathbb{1}_{X_i > t} \log \left(\frac{X_i}{t} \right),$$
$$n(t) = \sum_{i=1}^n \mathbb{1}_{X_i > t}.$$

Theorem 1.8

Suppose that \bar{F} has extreme value index $\gamma > 0$, and choose the threshold $t = t_n$ such that $t_n \ll n^{\frac{1}{\gamma}}$. Then

$$\hat{\gamma}_n = \hat{\gamma}_n(t_n) \xrightarrow{P} \gamma, \quad n \rightarrow \infty.$$

Proof of Theorem 1.8 (1/3)

- For any $\epsilon > 0$, we have

$$\begin{aligned} & P(|\hat{\gamma}_n - \gamma| > \epsilon) \\ &= \sum_{k=0}^{\infty} P(|\hat{\gamma}_n - \gamma| > \epsilon \mid n(t) = k) \cdot P(n(t) = k) \end{aligned}$$

- It thus suffices to show that

- ① $n(t_n) \xrightarrow{P} \infty$ as $n \rightarrow \infty$, and
- ② conditionally on $n(t) = k$, $\hat{\gamma}_n \xrightarrow{P} \gamma$ as $k \rightarrow \infty$.

Proof of (1):

Proof of Theorem 1.8 (2/3)

$$\hat{\gamma}_n = \hat{\gamma}_n(t) = \frac{1}{n(t)} \sum_{i=1}^n \mathbb{1}_{X_i > t} \log \left(\frac{X_i}{t} \right)$$

- ① $n(t_n) \xrightarrow{P} \infty$ as $n \rightarrow \infty$, and
- ② Conditionally on $n(t) = k$, $\hat{\gamma}_n \xrightarrow{P} \gamma$ as $k \rightarrow \infty$.

Proof of (2):

Proof of Theorem 1.8 (3/3)

$$\begin{aligned} E\left(\log \frac{X}{t} \mid X > t\right) &= \int_0^\infty \frac{\overline{F}(t \exp(s))}{\overline{F}(t)} ds \\ &\rightarrow \int_0^\infty \exp(s)^{-1/\gamma} ds = \gamma \end{aligned}$$

Theorem 1.9 (Karamata's representation, dHF B.1.6)

Let $f : [0, \infty) \rightarrow \mathbb{R}$ be regularly varying with index α , then there exist functions $a(t)$ and $c(t)$ with

$$c(t) \xrightarrow{t \rightarrow \infty} c_0 \in (0, \infty), \quad a(t) \xrightarrow{t \rightarrow \infty} \alpha$$

and a $t_0 > 0$ such that

$$f(t) = c(t) \exp\left(\int_{t_0}^t \frac{a(s)}{s} ds\right), \quad t > t_0.$$

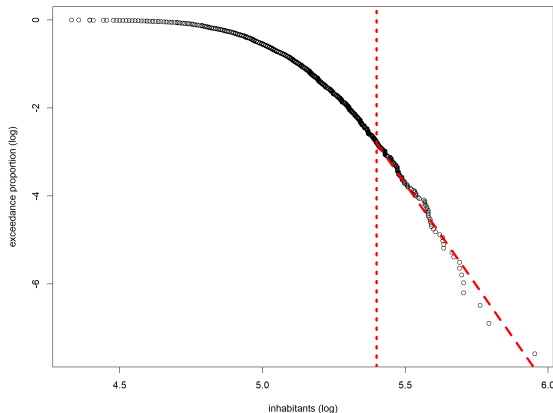
Corollary 1.10

Let f be regularly varying with index α . Then for any $\epsilon > 0$, there exists some $C = C(\epsilon)$ such that

$$f(t) \leq C(\epsilon) t^{\alpha + \epsilon}.$$

Application to the water levels at Hoek van Holland

Log-log plot of empirical survival function



- ▶ if F has extreme value index $\gamma > 0$, then

$$\bar{F}(tx) \approx \bar{F}(t)x^{-1/\gamma}, \quad x \rightarrow \infty$$

- ▶ Estimate the extreme value index by $\hat{\gamma}_n$ as above, with some suitable threshold $t = t_n$.
- ▶ Extrapolate the tail of the distribution as

$$\tilde{\bar{F}}(y) = \hat{\bar{F}}(t) \left(\frac{y}{t} \right)^{-1/\hat{\gamma}_n}$$

Exercise

Find a corresponding estimate for a large quantile $F^{-1}(1 - \alpha)$, for $\alpha \approx 0$.

Exercise (optional)

- 1 Choose a continuous distribution F that has a regularly varying tail with extreme value index $\gamma > 0$. In particular, prove this regular variation for your choice of F .
- 2 Simulate data from this distribution.
- 3 Pretend that you don't know the real value of γ , and use $\hat{\gamma}_n$ to estimate γ . Choose a suitable value of t . Compute the error $e = \hat{\gamma}^H - \gamma$.
- 4 Repeat steps (2) and (3) for 500 times. You collect $\{e_1, \dots, e_{500}\}$. Make a boxplot and histogram on the errors. Comment.