

Exercise session 1: univariate extremes

Anna Kiriliouk and Johan Segers

UCLouvain, October 2024

You'll need to install and load the `ismev` package first using `install.packages("ismev")` and `library(ismev)`.

Exercise 1.

Load the `ismev` package and the data `Ex1uni.RData` as below. The `data` vector contains daily maximal speeds of wind gusts, measured in kilometers per hour in Eindhoven (the Netherlands) during extended winter (October–March), from October 2001 up to and including March 2022. The `times` vector contains the dates of the observations.

```
library(ismev)
load("Ex1uni.RData")
```

- a. Use the following command to select monthly maxima. Also select the yearly maxima (careful: one for each extended winter season, spanning from October to March).

```
monthly <- as.vector(tapply(data, paste(times$year,times$mon), max))

yearly <- apply(matrix(monthly,ncol = 21, nrow = 6), 2, max)
```

- b. Plot the two datasets. Do they look stationary?

```
plot(monthly, ylim = c(35,130))

plot(yearly, ylim = c(35,130))

# The data looks stationary.
```

- c. Fit a GEV distribution to both the monthly and the yearly maxima and check the goodness-of-fit plots. You'll need to use the functions `gev.fit` and `gev.diag`. Would you choose to continue with the monthly or the yearly maxima?

```
gmonthly <- gev.fit(monthly, show = FALSE)
gyearly <- gev.fit(yearly, show = FALSE)
gev.diag(gmonthly)

gev.diag(gyearly)
```

I would go with the monthly maxima: the fit seems satisfactory and confidence intervals would be narrow

- d. For the chosen data (monthly or weekly), give the GEV parameter estimates with 95 % confidence intervals. Is the data bounded-tailed, light-tailed or heavy-tailed?

```
round(cbind(gmonthly$mle - 1.96*gmonthly$se, gmonthly$mle, gmonthly$mle + 1.96*gmonthly$se),2)

##      [,1] [,2] [,3]
## [1,] 62.71 65.26 67.82
## [2,] 11.18 13.01 14.84
```

```
## [3,] -0.18 -0.06 0.07
round(cbind(gyearly$mle - 1.96*gyearly$se, gyearly$mle, gyearly$mle + 1.96*gyearly$se),2)

##      [,1] [,2] [,3]
## [1,] 82.17 88.15 94.14
## [2,]  8.06 12.32 16.58
## [3,] -0.49 -0.15  0.20
# Light-tailed, in the Gumbel domain of attraction
```

- e. Using the function `gum.fit`, estimate the parameters of a GEV distribution with $\xi = 0$. Decide whether such a model suffices based on a likelihood ratio test. *Reminder: if $L(\xi, \sigma, \mu)$ denotes the log-likelihood function, then $2\{L(\xi, \sigma, \mu) - L(0, \sigma, \mu)\} \rightarrow \chi_1^2$.*

```
gummonthly <- gum.fit(monthly, show = FALSE)
gumyearly <- gum.fit(yearly, show = FALSE)
gum.diag(gummonthly)

gum.diag(gumyearly)

(Dmonthly <- 2*(gummonthly$nullh - gmonthly$nullh))

## [1] 0.7391278
(Dyearly <- 2*(gumyearly$nullh - gyearly$nullh))
```

```
## [1] 0.6194026
qchisq(0.95, df = 1)

## [1] 3.841459
```

- f. The return-level plot returned by `gum.diag` is not very clear. Using the formula seen this morning (slide 20), calculate the 10-year and the 1000-year return levels. Confidence intervals can be obtained using the delta method: if $\xi = 0$, $\nabla x_p^T = (1, -\log(-\log(1-p)))$, and the covariance matrix of $(\hat{\mu}, \hat{\sigma})$ can be obtained from the output of `gum.fit`. Is there a big difference between the 10-year and the 1000-year return levels? Is this surprising?

```
retLevGum <- function(pars, cov, p){
  rl <- pars[1] - pars[2]*log(-log(1-p))
  grd <- c(1, -log(-log(1-p)))
  se <- sqrt(t(grd) %*% cov %*% grd)
  return(c(rl - 1.96*se, rl, rl + 1.96*se))
}
retLevGum(gummonthly$mle, gummonthly$cov, (1/10)/6) #convert to yearly

## [1] 109.0606 117.1775 125.2943
retLevGum(gummonthly$mle, gummonthly$cov, (1/1000)/6) #convert to yearly

## [1] 160.3162 176.2404 192.1645
retLevGum(gumyearly$mle, gumyearly$cov, 1/10)

## [1] 102.2971 113.9311 125.5652
retLevGum(gumyearly$mle, gumyearly$cov, 1/1000)

## [1] 140.2321 169.2167 198.2013
```

- g. Fit a GP distribution to the daily maximal speeds of wind gusts that exceed 60 km/h and check the goodness-of-fit plots. You'll need to use the functions `gpd.fit` and `gpd.diag`.

```
gp60 <- gpd.fit(data, threshold = 60, show = FALSE)
gpd.diag(gp60)
```

- h. Use the functions `gpd.fitrange` and `mrl.plot` to decide whether the 60 km/h threshold is appropriate: if not, what threshold would you suggest?

Next, estimate the probability of a wind gust exceeding 144 km/h (the harshest inland wind ever to be measured in the Netherlands, during storm Eunice in February 2022) using the semi-parametric model on slide 36.

```
gpd.fitrange(data, umin = 50, umax = 80, nint = 10, show = FALSE)
```

```
# hard to say based on above plots: slightly higher (70) seems better
mrl.plot(data, umin = 30, umax = 100)
```

```
# impossible to say something based on the above plot
zeta_u <- length(which(data > 70))/length(data)
gpres <- gpd.fit(data, threshold = 70, show = FALSE)
xi <- gpres$mle[2]
alpha <- gpres$mle[1]
(prob <- zeta_u*(1 + xi*(144-70)/alpha)^(-1/xi))
```

```
## [1] 1.207147e-06
```

Exercise 2.

Load the `isnev` package and the data `fremantle` from the same package. This gives a `data.frame` where the first column contains the years, the second column gives annual maximum sea levels recorded at Fremantle, Western Australia, and the third column gives annual mean values of the Southern Oscillation Index, which is a proxy for meteorological volatility.

```
library(isnev)
data(fremantle)
```

- a. Plot the data. Would a stationary GEV model be appropriate?

```
plot(fremantle$Year, fremantle$SeaLevel, ylim = c(1,2))
```

```
# some evidence of non-stationarity
```

- b. Fit three non-stationary GEV models: for fixed scale and shape, try letting the location

- vary linearly with time
- vary linearly with the southern oscillation index
- vary linearly with both time and the southern oscillation index

Which model would you choose?

```
covar <- matrix(ncol = 2, nrow = 86)
covar[, 2] <- fremantle[, 3]
covar[, 1] <- c(1:86)
M0 <- gev.fit(fremantle$SeaLevel, show = FALSE)
M1 <- gev.fit(fremantle$SeaLevel, ydat = covar, mul = 1, show = FALSE)
M2 <- gev.fit(fremantle$SeaLevel, ydat = covar, mul = 2, show = FALSE)
M3 <- gev.fit(fremantle$SeaLevel, ydat = covar, mul = c(1,2), show = FALSE)
```

```
# AICs:
2*3+ 2*M0$nlh

## [1] -81.13326

2*4+ 2*M1$nlh

## [1] -91.57944

2*4+ 2*M2$nlh

## [1] -86.42228

2*5+ 2*M3$nlh

## [1] -97.65139
```

Exercise 3.

Simulate samples of size $n \in \{2000, 20000, 100000\}$ from the following two distributions:

- The inverse gamma distribution with shape $\alpha = 3$ and rate $\beta = 1$ (you can take the reciprocal of a gamma random variable with the same shape and rate).
- The log-Pareto distribution, whose cumulative distribution function is $F(x) = 1/\log(x)$ for $x \geq e$.

Using blocks of size 100, calculate the block maxima for each of the samples (leading to 20, 200, and 1000 block maxima respectively) and fit a GEV distribution (if possible!). Are these two distributions in the max-domain of attraction of a GEV? If yes, with what shape parameter?

```
library(ismev)
set.seed(1)
sample <- 1/rgamma(n = 100000, shape = 3, rate = 1)
maxima <- apply(matrix(sample, nrow = 100), 2, max)
gfit1 <- gev.fit(maxima[1:20], show = FALSE)
gfit2 <- gev.fit(maxima[1:200], show = FALSE)
gfit3 <- gev.fit(maxima, show = FALSE)
gev.diag(gfit1)

gev.diag(gfit2)

gev.diag(gfit3)

gfit1$mle

## [1] 2.2337537 0.5849453 0.4462718

gfit2$mle

## [1] 2.191086 0.820968 0.413875

gfit3$mle

## [1] 2.2780350 0.8601010 0.3348836

set.seed(1)
sample <- exp(1/runif(n = 100000, min = 0, max= 1))
summary(sample)

##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##      2.718   3.788   7.391    Inf  56.447    Inf
```

too heavy-tailed!