

An evaluation of automated GPD threshold selection methods for hydrological extremes across different scales

Stelian Curceac^{a,*}, Peter M. Atkinson^{b,c,d,e}, Alice Milne^f, Lianhai Wu^a, Paul Harris^a

^a Rothamsted Research, Department of Sustainable Agriculture Sciences, North Wyke EX20 2SB, Devon, UK.

^b Lancaster Environment Centre, Lancaster University, Lancaster LA1 4YQ, UK.

^c Geography and Environment, University of Southampton, Highfield, Southampton SO17 1BJ, UK.

^d School of Geography, Archaeology and Palaeoecology, Queen's University Belfast, BT7 1NN, Northern Ireland, UK.

^e State Key Laboratory of Resources and Environmental Information System, Institute of Geographical Sciences and Natural Resources Research, Chinese Academy of Sciences, Beijing 100101, China.

^f Rothamsted Research, Department of Sustainable Agriculture Sciences, Harpenden AL5 2JQ, UK

*Corresponding author: stelian.curceac@rothamsted.ac.uk, (+44) 7375596777

Resubmitted to: Journal of Hydrology

21/01/2020

Abstract

26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44

This study investigated core components of an extreme value methodology for the estimation of high-flow frequencies from agricultural surface water run-off. The Generalized Pareto distribution (GPD) was used to model excesses in time-series data that resulted from the 'Peaks Over Threshold' (POT) method. First, the performance of eight different GPD parameter estimators was evaluated through a Monte Carlo experiment. Second, building on the estimator comparison, two existing automated GPD threshold selection methods were evaluated against a proposed approach that automates the threshold stability plots. For this second experiment, methods were applied to discharge measured at a highly-instrumented agricultural research facility in the UK. By averaging fine-resolution 15-minute data to hourly, 6-hourly and daily scales, we were also able to determine the effect of scale on threshold selection, as well as the performance of each method. The results demonstrate the advantages of the proposed threshold selection method over two commonly applied methods, while at the same time providing useful insights into the effect of the choice of the scale of measurement on threshold selection. The results can be generalized to similar water monitoring schemes and are important for improved characterizations of flood events and the design of associated disaster management protocols.

Keywords: Generalized Pareto Distribution; Peaks over threshold; Threshold selection; Flood Frequency Analysis; Scale effects; Grassland agriculture.

45 1. Introduction

46 The magnitude and frequency of floods is likely to increase as a result of climate change (Bates
47 et al., 2008; Field et al., 2012; Kundzewicz et al., 2007) and this could push ecosystems beyond
48 the threshold of normal disturbance resulting in negative impacts that may be irreversible
49 (e.g. Thibault & Brown, 2008). Floods increase surface run-off, intensify erosion and introduce
50 more soil, organic matter and pollutants into water courses. Floods in areas of steep and
51 unstable slopes increase the possibility of landslides (Clarke & Rendell, 2006). Moreover,
52 increased runoff and flooding generally result in higher sediments and nutrient losses that
53 can lead to soil degradation (Bouraoui et al., 2004). They can have severe impacts on key
54 ecosystem services, such as those of support (e.g. water, nutrient cycling and soil protection),
55 regulation (e.g. climate) and culture (e.g. scenic recreation) (MA, 2005).

56 Flood Frequency Analysis (FFA) is a classic method to analyze the relationship between flood
57 magnitude and the corresponding frequency of occurrence. Reliable estimation and
58 prediction of high flow quantiles require extrapolation beyond the observed range of events,
59 commonly using parametric probability distributions. There are two main approaches for
60 defining extreme events in stationary time-series. The first is the block (usually annual)
61 maxima (AM) method where the dataset is divided into contiguous blocks of equal size and
62 the maximum values in each segment are considered. According to the Fisher-Tippett theorem
63 (Fisher & Tippett, 1928), these identically, independently distributed (iid) random variables
64 asymptotically follow a Generalized Extreme Value (GEV) distribution (Coles, 2001; Jenkinson,
65 1955). The second approach is known as the peaks-over threshold (POT) method, which
66 considers the values X that exceed a fixed high threshold u . The distribution function of the
67 excess values $X - u$, conditional on $X > u$, is a Generalized Pareto Distribution (GPD)

68 (Pickands, 1975). The case study we consider, contains six years of fine resolution (15-minute)
69 flow measurements, which is insufficient for effective fitting of the GEV distribution.
70 Therefore, only the POT method with the GDP was investigated.

71 The above two families of distributions have fundamental differences, but also theoretical
72 links (see Langousis et al., 2016). The GEV distribution is usually best fitted to annual maxima
73 samples and for this reason long historic records are required. This restriction does not apply
74 to the POT method since it includes all the peaks above a certain threshold allowing for
75 greater flexibility. The threshold must be large enough for the excesses to follow a GPD, but
76 an over-estimated threshold leads to reduced sample size and increases the variance of the
77 estimates. A smaller threshold increases the sample size but also the bias of the estimates as
78 the empirical distribution deviates from a perfect GPD model (Scarrott and MacDonald, 2012).
79 Clearly, GPD threshold selection is of key importance and there is no universally recognized
80 best performing method although various techniques have been proposed (see e.g. Langousis
81 et al. 2016 and Scarrott & MacDonald, 2012). Among them are probabilistic-based techniques
82 (Beirlant et al., 1996, 2006; Choulakian & Stephens, 2001; Deidda & Puliga, 2006; Goegebeur
83 et al., 2008; Hill, 1975), computational approaches (Beirlant et al., 2005; Danielsson et al.
84 2001; Hall, 1990; Thompson et al., 2009; Zoglat et al., 2014) and mixture models (Behrens et
85 al., 2004; Eastoe & Tawn, 2010; Solari & Losada, 2012). Graphical methods (Das & Ghosh,
86 2013; Deidda, 2010; Lang et al., 1999; Tanaka & Takara, 2010), such as the Mean Residual Life
87 (MRL) plot (Coles 2001; Beguería, 2005; Davison & Smith, 1990) are used commonly for the
88 selection of an optimal threshold, but have been criticized for the difficulty and subjectivity
89 of their interpretation (Scarrott & MacDonald 2012; Yang et al., 2018). Alternatively,
90 analytical methods have the advantage that they can be automated, and the associated

91 uncertainty can be quantified. Solari et al. (2017) proposed an automated threshold selection
92 method based on AD goodness of fit test. The application of their technique on long records
93 of precipitation and flow resulted in estimated thresholds that were within the stability
94 regions of the shape and modified scale parameters. Durocher et al. (2018) compared several
95 automatic methods and proposed a hybrid one where consistency with shape stability was
96 found for most of the considered sites.

97 In this study, we propose an empirical automated method for threshold determination, based
98 on threshold stability, which is evaluated against two commonly applied analytical methods,
99 together with eight alternatives for GDP parameter estimation. Furthermore, by averaging
100 the case study's 15-minute flow data to hourly, 6-hourly and daily supports, we determine
101 the effects of temporal measurement scale on threshold selection, as well as the performance
102 of each method.

103 The remainder of this paper is organized as follows. Section 2 presents the methods for GDP
104 parameter estimation, two analytical threshold selection techniques, this study's proposed
105 automated threshold stability method, and model evaluation diagnostics and indices. Section
106 3 describes the case study site and flow data, together with the simulation experiment design
107 used to evaluate the performance of the different GDP parameter estimators. Results are
108 presented in Section 4, which includes an investigation of scale effects through a series of
109 flow data integrations. Sections 5 and 6 discuss and conclude the study, respectively.

110 2. Methodology

111 The cumulative distribution function (CDF) of the iid excesses over an appropriate threshold
112 u for the GPD is:

$$G(x) = \Pr(X - u < x | X > u) = \begin{cases} 1 - \left(1 + \frac{\xi(x-u)}{\sigma}\right)^{-\frac{1}{\xi}}, & \xi \neq 0 \\ 1 - e^{\left(-\frac{x-u}{\sigma}\right)}, & \xi = 0 \end{cases}$$

where x , for this study, is the extreme flow in m^3s^{-1} , u is the location parameter, σ is the scale parameter and ξ is the shape parameter. The value of the shape parameter defines the type of distribution from the GPD family, that is, $\xi = 0$ refers to the exponential distribution, for $\xi > 0$ the corresponding distribution has a heavy upper tail that behaves like a power function with exponent $-1/\xi$ and for $\xi = 1$ the distribution is uniform. The Pareto distribution is obtained when $\xi < 0$.

2.1 GPD parameter estimators

The excesses above a suitable threshold are modelled by the GPD and the parameters of the distribution can be estimated by competing methods, where the Maximum Likelihood estimator (MLE) is the most commonly used (Prescott & Walden, 1980, 1983; Smith, 1985). Hosking and Wallis (1987) showed that MLE provides greater variance and bias for small samples compared to the Probability Weighted Moment (PWM) (Greenwood et al., 1979; Landwehr et al., 1979) and the Method of Moments (MOM) estimators. Coles and Dixon (1999) proposed a modified MLE which contains a penalty function for the shape parameter (i.e. the Maximum Penalized Likelihood estimator (MPLE). Zhang (2007) presented a hybrid Likelihood Moment estimator (LME) which provides feasible estimates and has high asymptotic efficiency. All of these methods are evaluated in this study, together with that suggested by Pickands (1975) and a maximum goodness-of-fit (MGF) estimator (e.g. Luceño, 2006). Estimator performance has been found to depend significantly on sample size and the value of the GPD shape parameter (Ashkar & Tatsambon, 2007; de Zea Bermudez & Kotz,

134 2010; Hosking & Wallis, 1987), and the choice of the estimator should be made based on the
135 specifics of the situation. The equations for the above estimators can be found in Appendix
136 A: Equations of the estimators.

137 2.2 Threshold selection methods

138 The selection of the threshold u is a crucial step in GDP extreme value analysis. On the one
139 hand, a small threshold results in a large sample that makes statistical inference more
140 effective, but can lead to biased estimates due to deviations of the empirical distributions
141 from the GPD model (e.g. Beirlant et al., 2005). On the other hand, when considering large
142 thresholds and consequently small samples, parameter estimates have a smaller expected
143 bias, but a larger variance that can be highly dependent on the estimation method. The two
144 main approaches for threshold selection are graphical methods, such as the MRL plot, and
145 analytical methods that can be automated.

146 An important assumption for the application of the POT method is that the extracted peaks
147 are independent. A commonly applied method is to use no more than 2-3 peaks per year
148 (Madsen et al., 1997; Todorovic, 1978) but it has been criticised for lack of flexibility. Another
149 solution is to consider a minimum separation interval between successive peaks (Cunnane,
150 1979; Lang et al., 1999). This minimum separation interval accords to the scale and nature of
151 the measured process, but for daily flow data, an interval of a few days commonly ensures
152 that the peaks are generated from different events (Engeland et al., 2004). The
153 autocorrelation function is a popular choice for the investigation of serial dependence in a
154 time series. However, this approach assumes normally distributed variables, which is not the
155 case for peak discharges, so other independence tests should be implemented (e.g. Ledford
156 and Tawn, 2003; Reiss and Thomas, 2007). In this study, and through prior experimentation,

157 maximum peaks separated by a minimum of three days were considered and their
158 independence was tested using Kendall's τ test (Claps and Laio, 2003; Ferguson et al., 2000).

159 2.2.1 Graphical methods: MRL plots

160 The most popular graphical method is the MRL plot (Coles, 2001; Davison & Smith, 1990). If
161 the scaled excesses $X_{u^*} = [X - u^* | X > u^*]$ above a threshold u^* are Generalized Pareto (GP)
162 distributed, then for every $u \geq u^*$, the scaled excesses $X_u = [X - u | X > u]$ are similarly GP
163 distributed with the same shape parameter ξ , a scale parameter $\sigma_u = \sigma_{u^*} + \xi(u - u^*)$ and
164 a mean value:

$$165 \bar{X}(u) = E[X - u | X > u] = \frac{\sigma_u}{1 - \xi} = \frac{\sigma_{u^*} + \xi(u - u^*)}{1 - \xi} = Au + B$$

166 where $A = \xi/(1 - \xi)$ and $B = (\sigma_{u^*} - \xi u^*)/(1 - \xi)$ are the respective slope and intercept
167 of the linear relation. The sample estimates of the mean excesses are then plotted for
168 different values of the threshold and the most appropriate is considered to be the one after
169 which the mean excesses follow a straight line (e.g. Das & Ghosh, 2013).

170 Another graphical technique is to plot the estimated shape and/or modified scale parameters
171 for different threshold candidates and select the one above which the estimates are constant
172 (Brodin & Rootzén, 2009; Bommier, 2014; Sigauke & Bere, 2017). The main criticism of
173 graphical methods is that the interpretation of the plot can be ambiguous or subjective as it
174 is usually unclear which part of the curve is linear (Scarrott & MacDonald, 2012). In this
175 respect, attempts have been made to automate (Langousis et al., 2016) and estimate the
176 uncertainty (Liang et al., 2019) of the graphical methods.

177 2.2.2 Analytical methods: Square Error and Normality of Differences

178 The Square Error (SE) method was developed by Zoglat et al. (2014) following the work of
179 Beirlant et al. (2005), and is implemented as follows. Let u_1, u_2, \dots, u_n be n equally spaced
180 increasing threshold candidates. For each of these thresholds, estimate the scale σ_{u_j} and
181 shape ξ_{u_j} parameters for $j = 1, \dots, n$. Find N_{u_j} the exceedances that correspond to each
182 threshold u_j and simulate m independent samples of size N_{u_j} from the GPD with parameters
183 σ_{u_j} and ξ_{u_j} . For each probability $a \in A = \{0.05, 0.1, \dots, 0.95\}$ and each $i = 1, \dots, m$ calculate
184 the quantiles q_{a,u_j}^i and compute $q_{a,u_j}^{sim} = \frac{1}{m} \sum_{i=1}^m q_{a,u_j}^i$. The optimal threshold is the one for
185 which the square error $SE_{u_j} = \sum_{a \in A} \left(q_{a,u_j}^{sim} / q_{a,u_j}^{obs} \right)^2$ between the simulated and the observed
186 quantiles is minimum. The selection of the threshold candidates u_j can be defined by the user
187 or as an automated process. For example, the smallest threshold can be set as zero or the
188 median and the maximum threshold set as a high percentile of the data.

189 An alternative analytical method for threshold selection was proposed by Thompson et al.
190 (2009). Again, let u_1, u_2, \dots, u_n be n equally spaced increasing threshold candidates. For the
191 excesses above the threshold u_j , $\hat{\sigma}_{u_j}$ and $\hat{\xi}_{u_j}$ are the MLEs of the scale and shape parameters,
192 respectively, for $j = 1, \dots, n$. If $u \leq u_{j-1} < u_j$ is an appropriate threshold then according to
193 Coles (2001), $\sigma_{u_{j-1}} = \sigma_u + \xi(u_{j-1} - u)$ and $\sigma_{u_j} = \sigma_u + \xi(u_j - u)$. Consequently, $\sigma_{u_j} -$
194 $\sigma_{u_{j-1}} = \xi(u_j - u_{j-1})$ and from standard maximum likelihood theory we have that $E[\hat{\sigma}_{u_j}] \approx$
195 σ_{u_j} and $E[\hat{\xi}_{u_j}] = \xi$ for any j such that $u_j > u$. Respectively, $E[\tau_{u_j} - \tau_{u_{j-1}}] \approx 0$, $j =$
196 $2, \dots, n$ for $\tau_{u_j} = \hat{\sigma}_{u_j} - \hat{\xi}_{u_j} u_j$, $j = 1, \dots, n$. It follows that $\tau_{u_j} - \tau_{u_{j-1}}$ approximately follows
197 a normal distribution. Thompson et al. (2009) suggest Pearson's Chi-square test to examine
198 the null hypothesis of normality. However, this test has been criticised for having inferior

199 power properties (Moore, 1986). For this reason, we also applied the Anderson-Darling,
 200 Cramer-von Mises, Kolmogorov-Smirnov and Shapiro-Francia normality tests (Thode, 2002).
 201 Regardless of which of the five normality tests are used, we refer to this method as the
 202 ‘Normality of Differences’ method. According to this approach, a suitable threshold $u \leq$
 203 $u_{j-1} < u_j$ is the one for which all the differences $\tau_{u_j} - \tau_{u_{j-1}}$ are approximately normally
 204 distributed. We selected the appropriate threshold as the one for which the p -value of $\tau_{u_j} -$
 205 $\tau_{u_{j-1}}, j = 2, \dots, n$ is above 0.05. A smaller threshold would be selected for a smaller p -value
 206 (e.g. 0.01).

207 2.2.3 Proposed method based on Threshold Stability

208 For this study, we propose an automated threshold selection method based on stability plots
 209 (Coles, 2001; Scarrott & MacDonald 2012). If the GPD is an appropriate model for the excesses
 210 above a threshold u , then for all larger thresholds $u^* > u$ it will also be suitable with the shape
 211 parameter being relatively constant. In other words, it is the approximately linear horizontal
 212 part on the shape parameters versus thresholds plot. This does not apply for the scale
 213 parameter σ_{u^*} , as it changes with the threshold $\sigma_{u^*} = \sigma_u + \xi(u^* - u)$. However, the
 214 modified scale parameter $\sigma_1 = \sigma_{u^*} - \xi u$ remains relatively constant. Therefore, we fit a cubic
 215 smoothing spline to this plot and calculate the rate of change at each of m consecutive steps.
 216 The cubic smoothing spline estimate \hat{f} of a function f in the model $Y_i = f(x_i) + \varepsilon_i$, is defined
 217 as the minimizer of $\sum_{i=1}^n \{Y_i - \hat{f}(x_i)\}^2 + \lambda \int \hat{f}''(x)^2 dx$, where λ is the smoothing parameter.
 218 The minimum change rate locates the part of the plot where the shape and the modified scale
 219 parameters reach a plateau.

220 A preliminary analysis showed that a smoothing parameter value of $\lambda = 0.4$ of the cubic spline
221 function was the most appropriate to avoid both over- and under-fitting. A total of $n =$
222 1000 threshold candidates were used in each case and a cubic spline was fitted to the
223 corresponding estimated shape and modified scale parameters. The numbers of the
224 consecutive steps for which the minimum change rate was calculated, were $m =$
225 25, 50, 75 and 100 which corresponds to 2.5%, 5%, 7.5% and 10%, respectively, of the total
226 number of fitted values, that is, the total threshold candidates n .

227 2.3 Evaluation procedure

228 Quantile-Quantile (Q-Q) plots are commonly used to investigate the efficiency of the
229 statistical inference of the fitted GPD models. To quantify the difference between the
230 theoretical and empirical quantiles for probabilities $\alpha \in A = \{0.95, 0.951, \dots, 0.999\}$, various
231 error and agreement diagnostics were calculated. Specifically, we calculated the Mean Square
232 Error (MSE) (e.g. Turan and Yurdusev, 2009), the Normalized Root Mean Square Error
233 (NRMSE) (e.g. Sheta and El-Sherif, 1999) and the Relative Index of Agreement ($RD \in [0,1]$)
234 (Krause et al., 2005; Willmott, 1981). For ideal model performance, both MSE and NRMSE
235 should tend to zero, while RD should tend to unity. The NRMSE was obtained by dividing the
236 root MSE with the difference between minimum and maximum values and, thus, was less
237 sensitive to very large values and provided a more robust diagnostic than MSE.

238 3. Study site and datasets

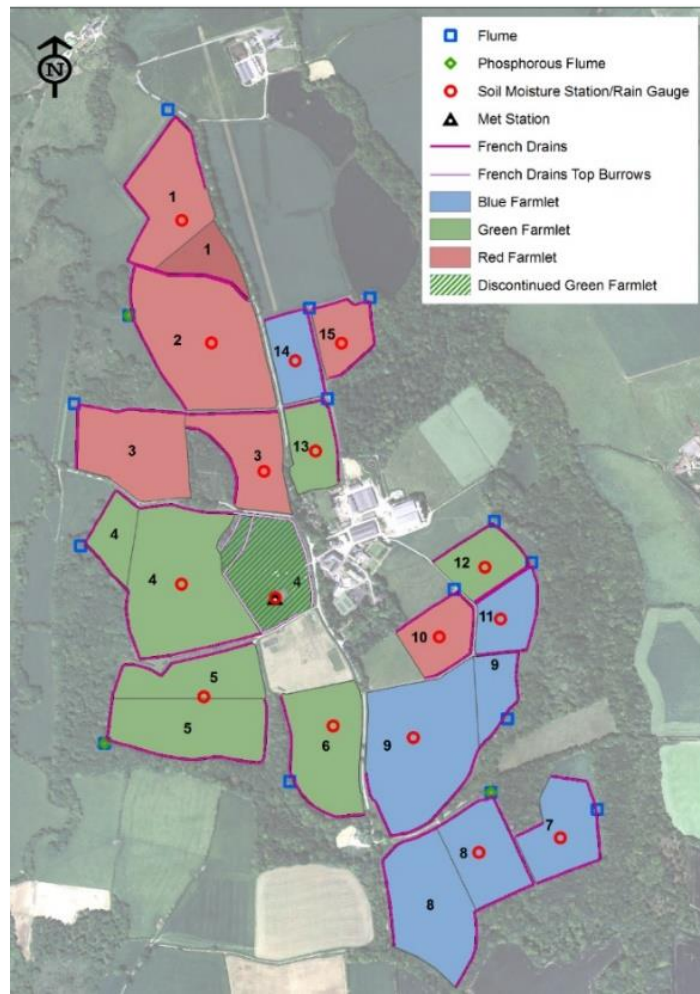
239 3.1 Study site

240 Flow discharge data come from a single sub-catchment of the North Wyke Farm Platform
241 (NWFP). The NWFP is a farm-scale experiment established in 2011 in the southwest of

242 England (50°46'10"N, 3°54'05"W) for research into sustainable grassland livestock systems
243 (Orr et al., 2016; Takahashi et al., 2018). The platform is located at an altitude in the range of
244 120-180 m above sea level. The platform's fields have a declining slope at the west towards
245 the River Taw and to the east, to one of its tributaries, the Cocktree stream. The soil texture
246 consists of a slightly stony clay loam topsoil (approximately 36% clay) above a mottled stony
247 clay (approximately 60% clay). The subsoil is impermeable to water and during rain events
248 most of the excess water moves by surface and sub-surface lateral flow towards the drainage
249 system described below.

250

251 Each of the 15 NWFP sub-catchments are hydrologically isolated through a combination of
252 topography and a network of French drains (800 mm deep trenches), which ensure that the
253 total runoff is channeled to instrumented flumes, measuring 15-minute water discharge and
254 water chemistry from October 2012. The discharge from each sub-catchment is measured
255 through a combination of primary and secondary flow devices (Liu et al., 2018). The primary
256 devices are H-type flumes (TRACOM Inc., Georgia, USA) with capacity designed for a 1-in-50
257 year storm event. The specific design of the H-type flume facilitates the accurate
258 measurement of both low and high flows and is relatively self-cleaning since it allows the
259 ready passage of sediment and particulate matter. A secondary flow measurement device
260 (OTT hydromet, Loveland, CO., USA) is used to measure the water height within the flume
261 and convert it to discharge rate using flume-specific formulas which depend on water height.
262 The flow is generated only from rainfall as the fields are not irrigated. At each sub-catchment,
263 15-minute precipitation and soil moisture are also monitored. (Figure 1).



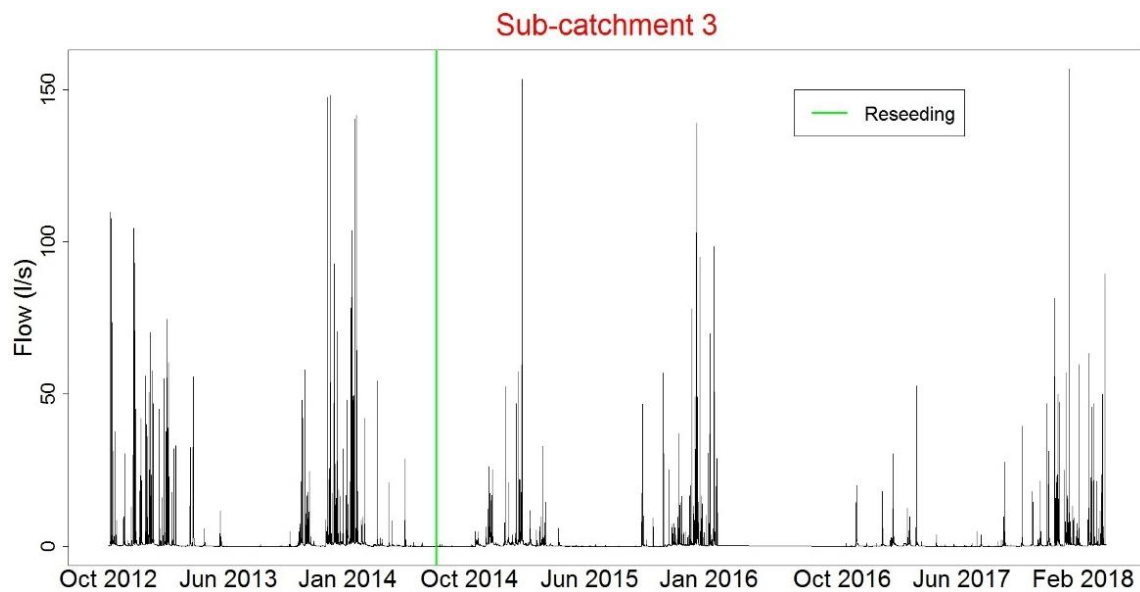
264

265 Figure 1: The three farmlets and the 15 sub-catchments of the North Wyke Farm Platform, with: (i)
 266 'blue' farmlet a mixture of white clover and high sugar perennial ryegrass; (ii) 'red' farmlet high
 267 sugar perennial ryegrass only and (iii) 'green' farmlet permanent pasture ("business as usual").

268 3.2 Measured data

269 For this study, we used the flow discharge measured at sub-catchment 3 of the NWFP, which
 270 is part of the 'red' farmlet (Figure 1) and 6.84 ha in size. Given this is a methodological-based
 271 study, we chose to use data from this sub-catchment as it has one of the smallest number of
 272 missing values (approximately 1%) for the six-year period (2012-2018). Imputation of the
 273 missing values was performed using a regularized iterative Principal Components Analysis
 274 (PCA impute) model (Josse & Husson, 2013). The largest imputed value was approximately 20

275 $l s^{-1}$ which is smaller than any threshold suggested (see below) and, therefore, is not
276 considered as a peak flow and does not affect the subsequent analysis. It should be noted
277 that, compared with measurements from many river or stream monitoring systems, the flow
278 data (Figure 2) are highly discontinuous with many zeros, as non-zero measurements occur
279 only after rainfall events.



280

281 Figure 2: Flow ($l s^{-1}$) measurements at sub-catchment 3 (2012 to 2018).

281

282 3.3 Simulated data

283 As a precursor to the empirical study, the performance of the eight GDP parameter estimators
284 was assessed through a Monte Carlo experiment. We generated random time-series of
285 different sample sizes ($n = 25, 50, 100, 250, 500, 1000$) from a GPD distribution with a
286 known shape parameter ($\xi = -0.5, -0.25, 0, 0.25$ and 0.5). For each combination, 10,000
287 random samples were generated. The performance of the estimators was evaluated using:
288 (a) bar plots for MSE values and (b) boxplots for estimated ξ . Here the “error” in MSE is the
289 difference between the actual (or known) ξ and that estimated, where MSE incorporates both

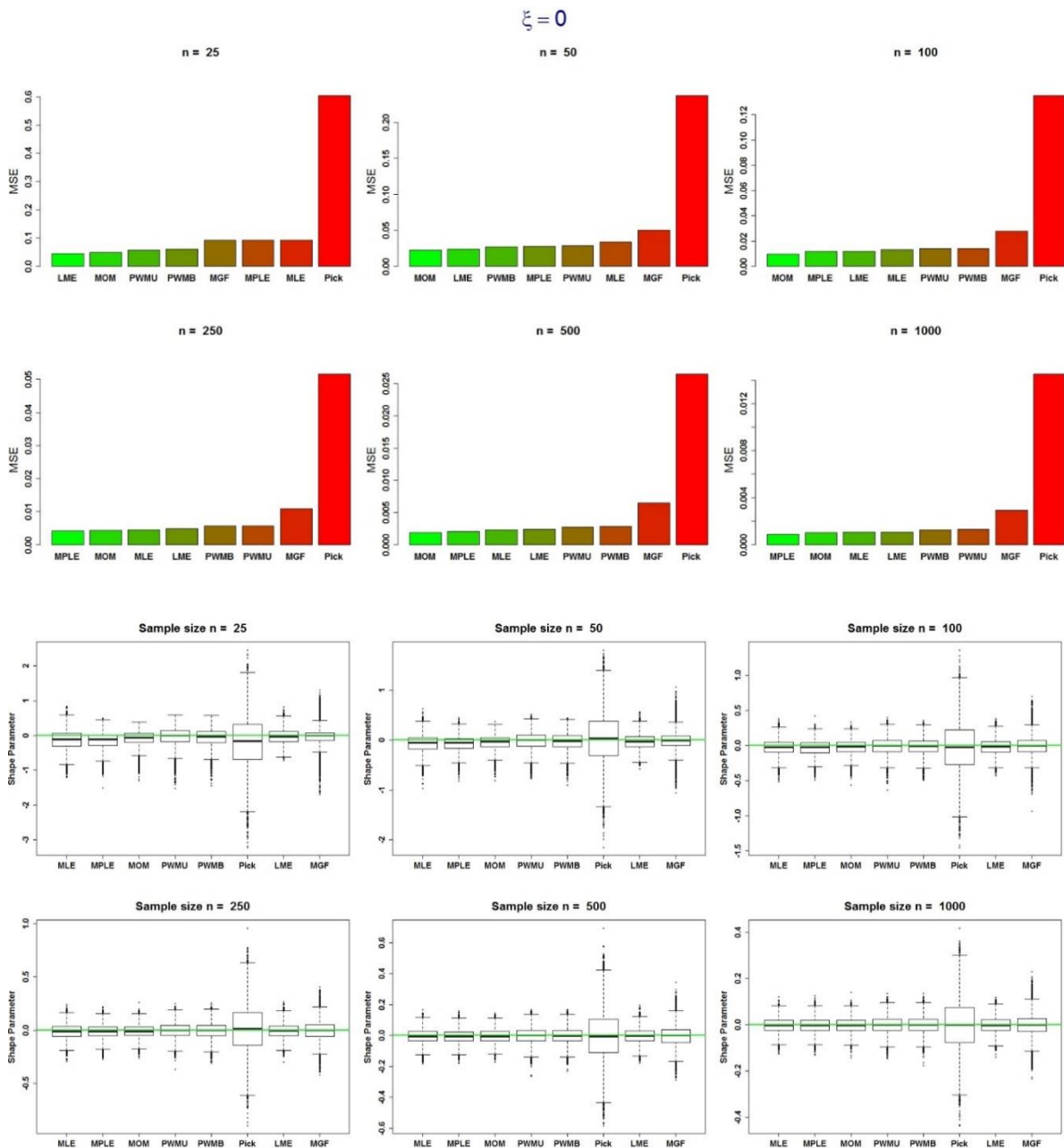
290 the variance and the bias of the estimators. Outcomes were used to guide the analyses with
291 the measured NWFP flow data.

292 4. Results

293 4.1 Monte Carlo study for Performance of GPD estimators

294 Our simulated data analysis showed that the performance of the GPD parameter estimators
295 depends on both the sample size n (see performance plots in Figure 3 for a shape parameter
296 of $\xi = 0$ only) and the value of the shape parameter ξ (see supplementary material for
297 performance plots with $\xi = -0.5, -0.25, 0.25$ and 0.5), which accords with previous studies (e.g.
298 Gharib et al., 2017; Mackay et al., 2011). On viewing all plots, the maximum likelihood (MLE
299 and MPLE) estimators were both negatively biased for small sample sizes for any value of the
300 shape parameter and their performance increased in terms of bias and variance as sample
301 size increased. The MLE outperformed the other estimators for large sample sizes for all
302 values of the shape parameter. The unbiased and biased probability weighted moments,
303 PWMU and PWMB respectively, were consistently the least biased amongst all estimators
304 and provided a small variance, which was less sensitive to sample size compared to the
305 likelihood estimators. According to the MSE, the PWM estimators were most appropriate for
306 small sample sizes and positive shape parameters. The MOM estimator had a similar behavior
307 to the PWMs when $\xi \leq 0$ but had a negative bias for $\xi > 0$ and the bias increased as the
308 value of the shape parameter and the sample size increased. Pickland's estimator ('Pick') and
309 the MGF estimators produced a large variance and the least accurate estimates of the shape
310 parameter, through the whole range of the examined values. LME was among the best
311 performing estimators regarding accuracy and bias, except for the very short tails ($\xi = 0.5$,

312 see supplementary material), when the estimates deviated greatly from the rest of the
 313 estimators and the predefined value of the shape parameter. In summary, the MLE/MPLE,
 314 PWMU/PWMB and the LME were considered the most unbiased and precise estimators and
 315 so we select only from this reduced group of estimators in subsequent analyses using the
 316 measured data.



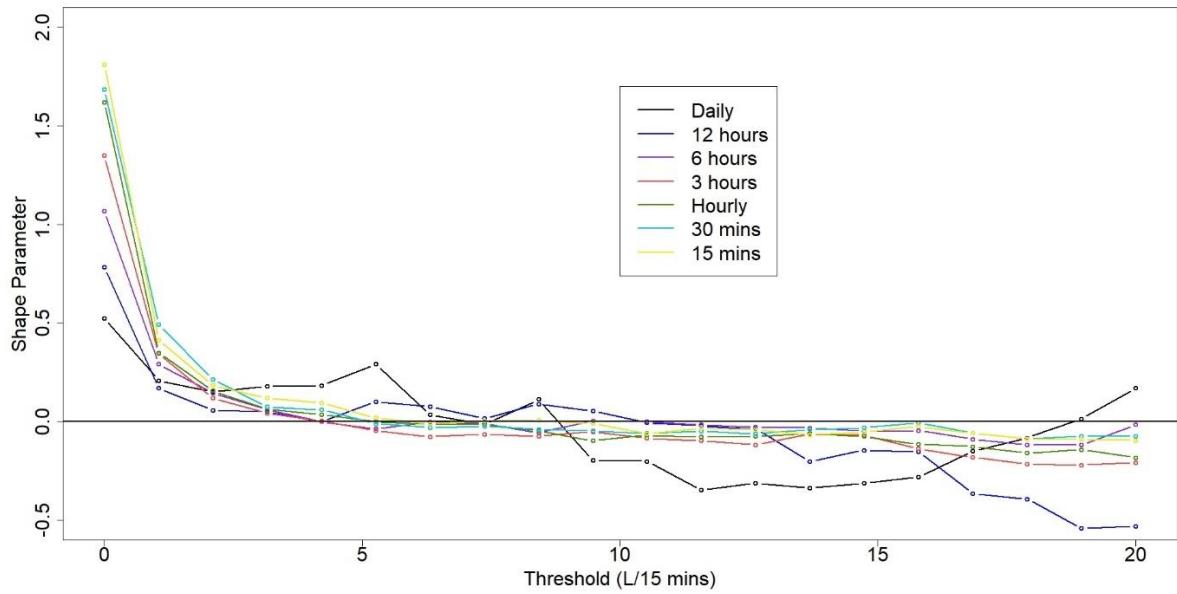
319 Figure 3: Performance of GPD estimators for shape parameter $\xi = 0$ and for six different sample sizes
 320 ($n = 25, 50, 100, 250, 500, 1000$).

321 4.2 Empirical study for Threshold Selection

322 4.2.1 Preliminary effects of data aggregation

323 Initially, the flow (l s^{-1}) time-series of 15-minute resolution was averaged to time-series data
324 of 30 minutes, hourly, 3-hourly, 6-hourly, 12-hourly and daily resolutions. Figure 4 shows the
325 behavior of the MLE-estimated shape parameters for a range of thresholds for the differently
326 aggregated flow data. The range of thresholds was set from the median to the maximum for
327 which daily flow can be fitted efficiently. The shape parameter is in the range of 0.5 to almost
328 2 for the minimum threshold, has a decreasing trend as the threshold increases and can
329 become negative for the largest thresholds. The similar shape characteristics could be an
330 indication that the shape parameter describes an inherent feature of the process and that
331 changes of scale, which affect the size or variability of the observed values of the process, do
332 not substantially change the shape characteristics of these observations. For the remainder
333 of this study, results from the 30-minute, 3-hourly and 12-hourly aggregations are not
334 reported as retained aggregations (hourly, 6-hourly and daily) communicate all key outcomes
335 adequately.

336 Kendall's τ test showed that the maximum peaks separated by a minimum of three days were
337 reasonably independent (Figure 5). The statistics τ are large for the lowest thresholds where
338 the peaks are numerous and autocorrelated. With an increasing threshold, the values of the
339 τ decrease rapidly and are below the 95% acceptance limits which supports the null
340 hypothesis of independence of the peaks.

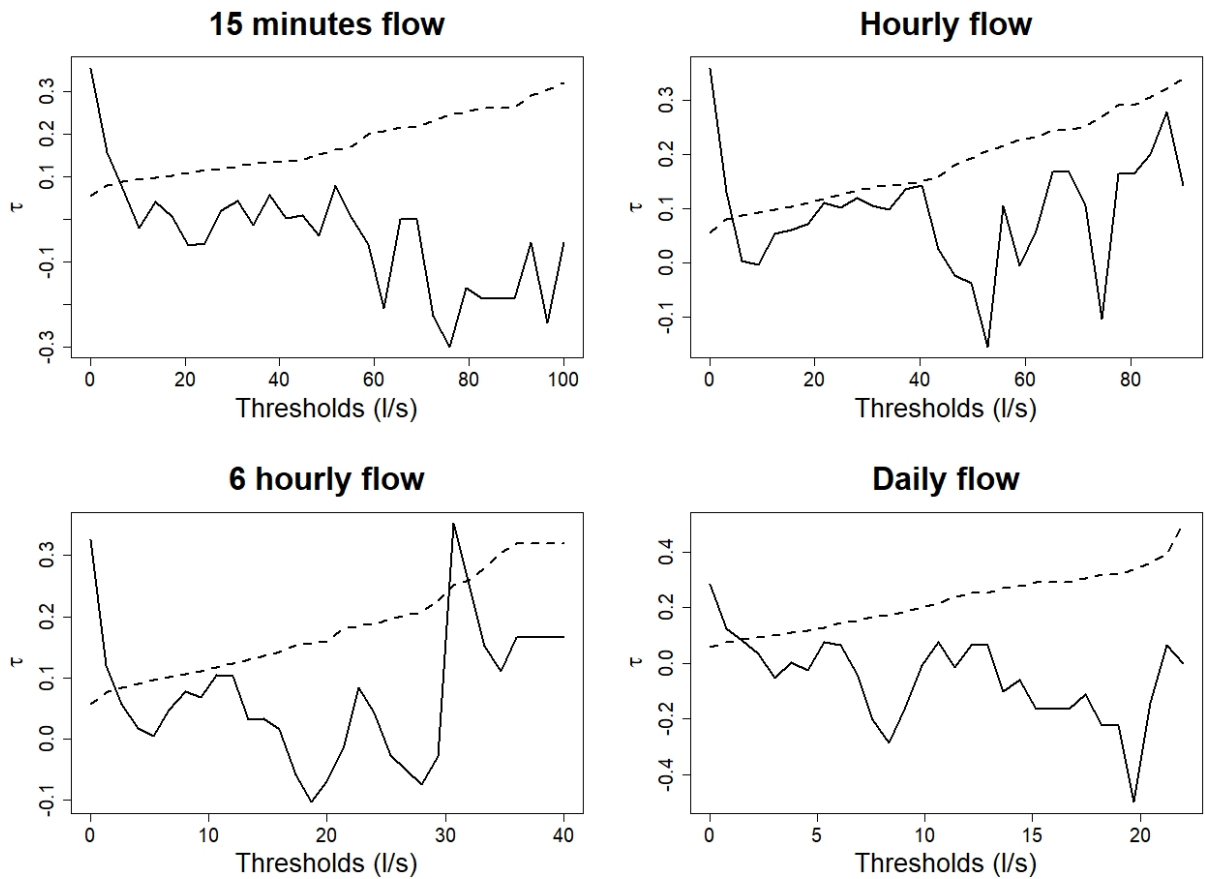


341

342

Figure 4: Shape parameter characteristics of measured (15-minute) and a series of averaged (30-minute to daily) flow rates.

343



344

345

Figure 5: Kendall's test statistic τ (solid lines) along with the 95% acceptance limits of the test (dashed lines).

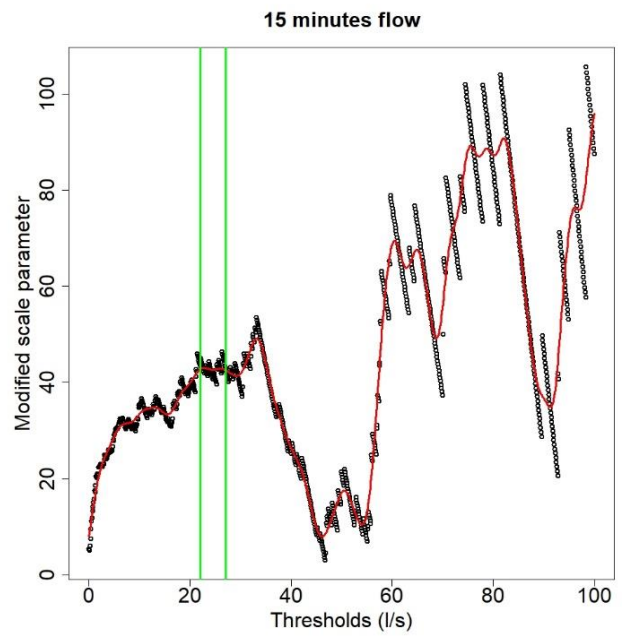
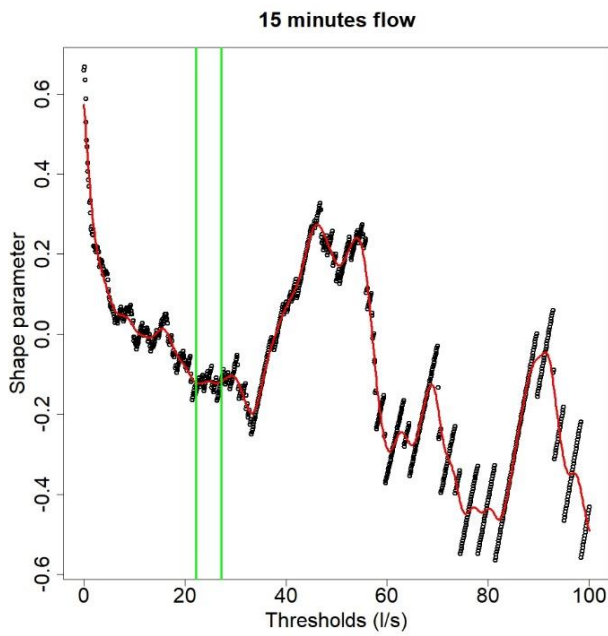
346

347 4.2.2 Automated Threshold Stability plots

348 The choice of estimators for the shape and modified scale parameters was guided by the
349 results of the Monte Carlo experiment (Section 4.1). For example, for thresholds $u_j = 1, 2, \dots, 5$
350 of the 15-minute flow data, the number of exceedances was $N_{u_j} > 300$ and the shape
351 parameter ξ_{u_j} between 0.5 and 0.25. For this combination, MLE, MPLE, PWMU, PWMB and
352 LME were the best performing estimators. Thus, for our empirical study, we choose LME due
353 to its consistently precise and unbiased estimates of positive shape parameters for a large
354 sample size. Increasing the thresholds u_j resulted in a reduced sample size ($100 < N_{u_j} <$
355 250) and negative values of the shape parameter. In this case, we choose MPLE for our
356 empirical work. In all the other cases, the PWMU estimator was preferred as it provided
357 unbiased estimates with small variance.

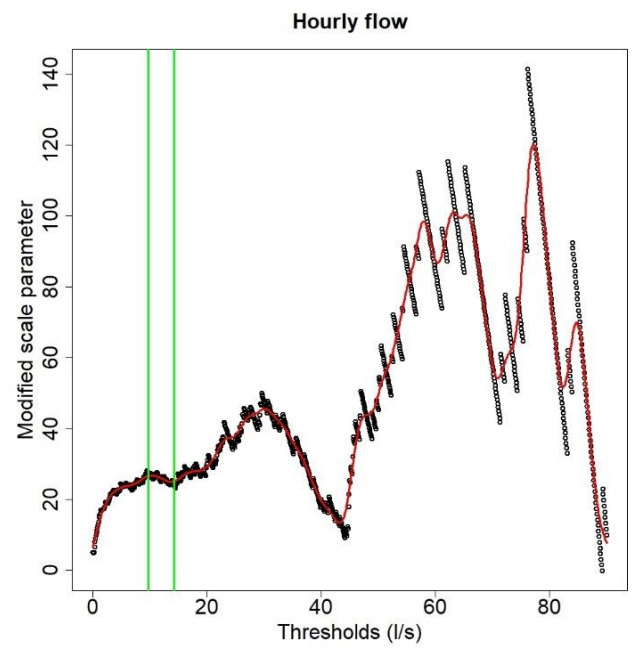
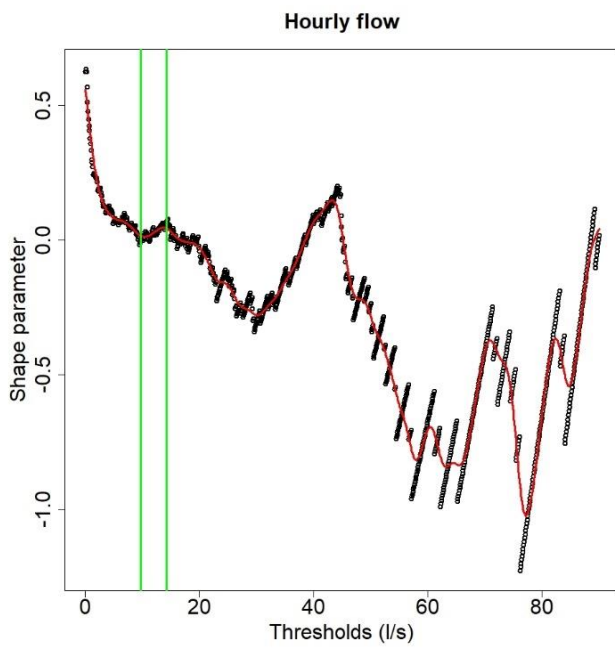
358 Stability plots are given in Figure 6 for different flow aggregations, where results reveal our
359 'Automated Threshold Stability' (ATS) extension to be reasonably robust, since changes in the
360 number of consecutive steps m had a very small impact on the selected threshold and usually
361 resulted in over-lapping regions from which the threshold was considered. The peak flows at
362 15 minutes and hourly resolution did not provide many regions that could be considered as a
363 plateau, so the number of consecutive steps was set to $m = 50$ (5% of the total) to also capture
364 the smaller approximately linear horizontal parts. Interestingly, for each aggregation, fitting
365 the same cubic spline functions to both the estimated shape and modified scale parameters,
366 resulted in almost identical suggested thresholds.

367 a)



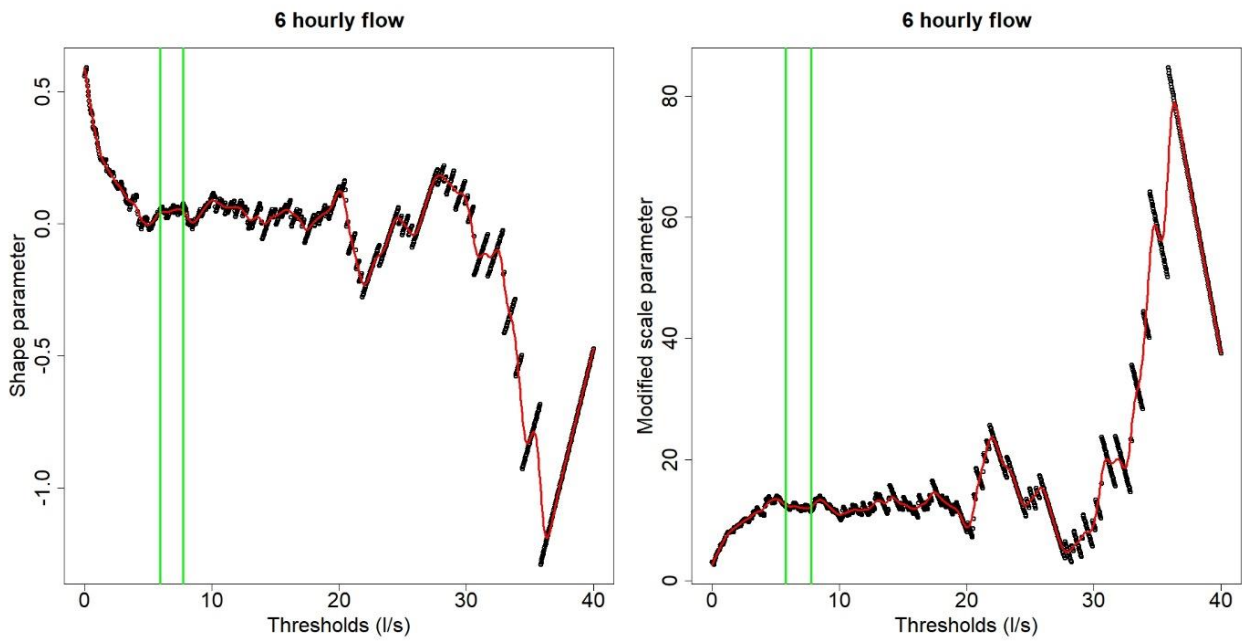
368

369 b)



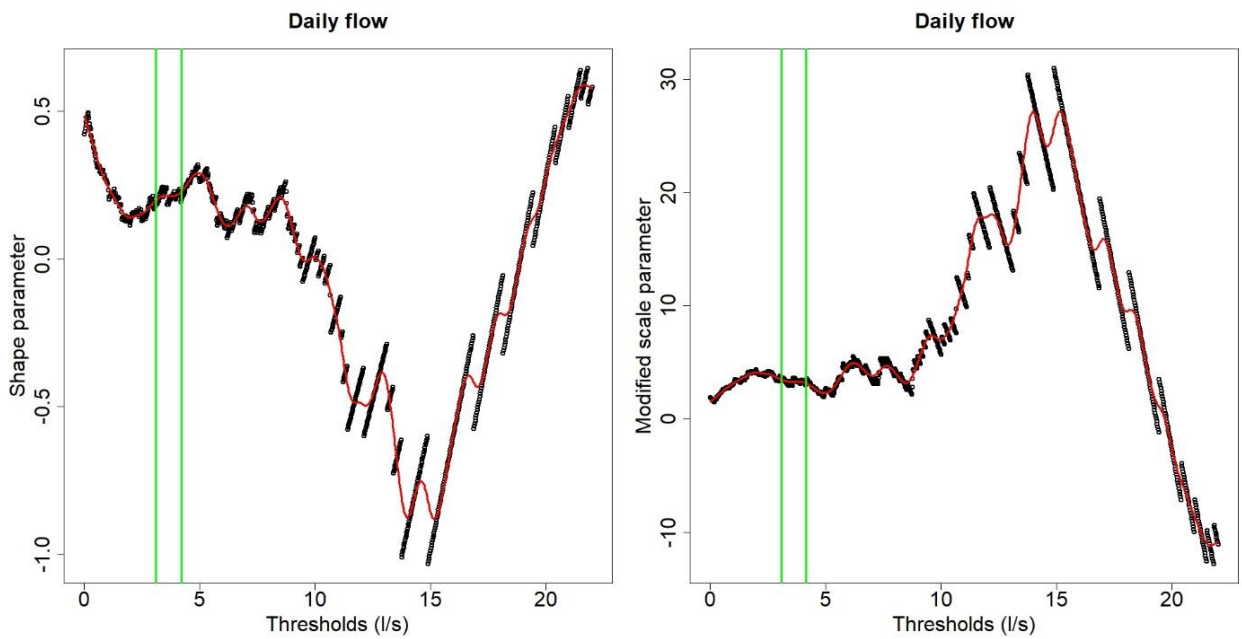
370

371 c)



372

373 d)



374

375 Figure 6: Automated Threshold Stability (ATS) method: Selected threshold (that between the vertical
376 green lines) of a) 15 minutes, b) hourly, c) 6 hourly and d) daily flow based on smoothing splines.

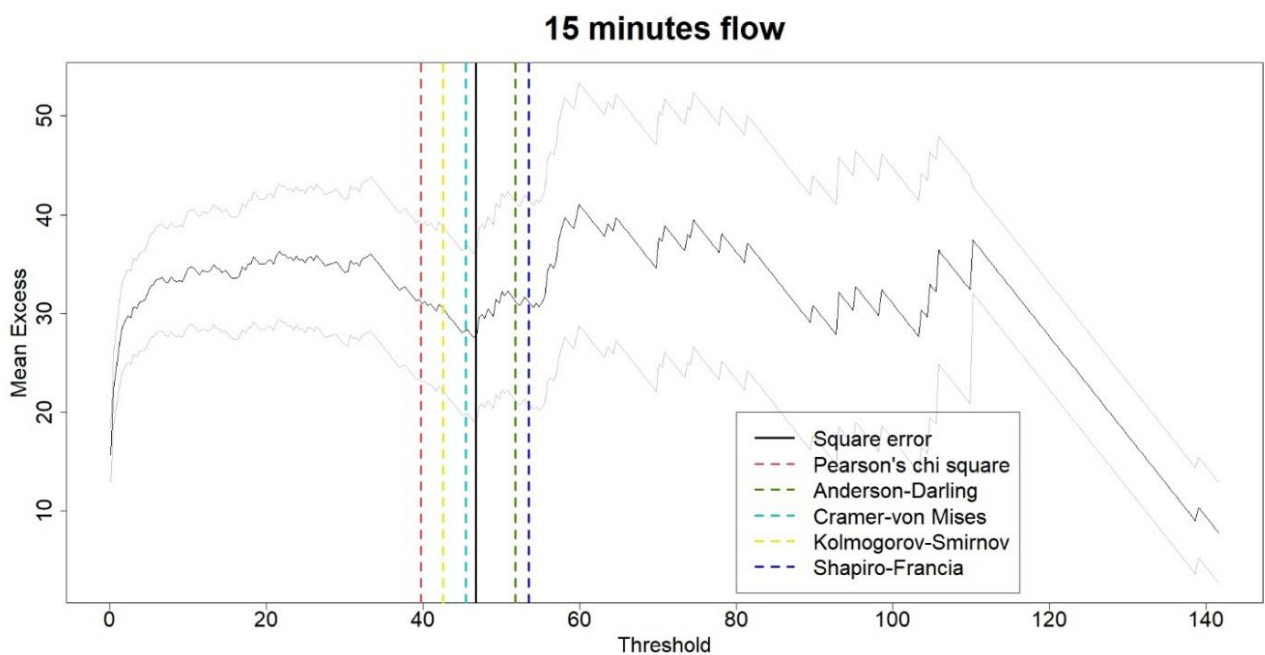
377 4.2.3 Analytical threshold selection methods: Square Error and Normality of Differences

378 The choice of GDP estimators for the simulation of the quantiles for the SE method was
379 performed using a similar procedure as described in Section 4.2.2, while the approach based
380 on the Normality of Differences test is based on assumptions of maximum likelihood theory,
381 and consequently the shape parameter was estimated by the MLE. The number n of the
382 considered thresholds u_n plays an important role in the results. Thompson et al. (2009)
383 suggested $n = 100$ and reported that for $n < 100$, less reliable results were obtained. We
384 similarly specified $n = 100$ but also found the thresholds to be over-estimated for $n > 100$.

385 Our results indicated little consistency in the selection of thresholds where a specific part of
386 the MRL plot could be considered approximately linear. The thresholds of the 15-minute peak
387 flow estimated by the SE method and the Normality of Differences tests (Figure 7a) are
388 considerably larger than that based on this study's ATS method (Figure 6a) at around 40 to 50
389 l/s and 20 to 30 l/s, respectively. Only for the daily flow data (Figure 7d), the threshold
390 estimated by the SE method was smaller than those estimated from the Normality of
391 Differences tests and relatively close to the threshold estimated by ATS (Figure 6d). For hourly
392 flow data (Figure 6b and Figure 7b), ATS and Pearson's chi square test (for Normality of
393 Differences) provided almost identical estimates, while all other methods suggested much
394 larger thresholds. Noticeably, the hourly thresholds estimated by the SE method and the
395 Shapiro-Francia test are very close at 44.68 l/s and 45.33 l/s, respectively (Figure 7b), but
396 result in considerably different shape parameters (Table 1). Figure 6b reveals hourly
397 thresholds to be in the region where the shape characteristics show large fluctuations due to
398 the small sample size that results in an inefficient fit of the GPD and likely spurious estimates
399 of the shape parameter.

400 The performance of the Normality of Differences method depended greatly on both the given
401 normality test and on data resolution. For the 15-minute flow data, all normality tests
402 provided relatively similar threshold selections (Figure 7a), which was not the case for the
403 hourly and 6-hourly flow data (Figure 7b and Figure 7c). For the daily flow data (Figure 7d),
404 thresholds were estimated too large and consequently result in too few values for efficient
405 statistical inference. In general, the smaller the selected threshold, given that the excesses
406 are satisfactorily modelled by the GPD, the lower the uncertainty and consequently the lower
407 the variance in the parameter estimates due to larger sample sizes.

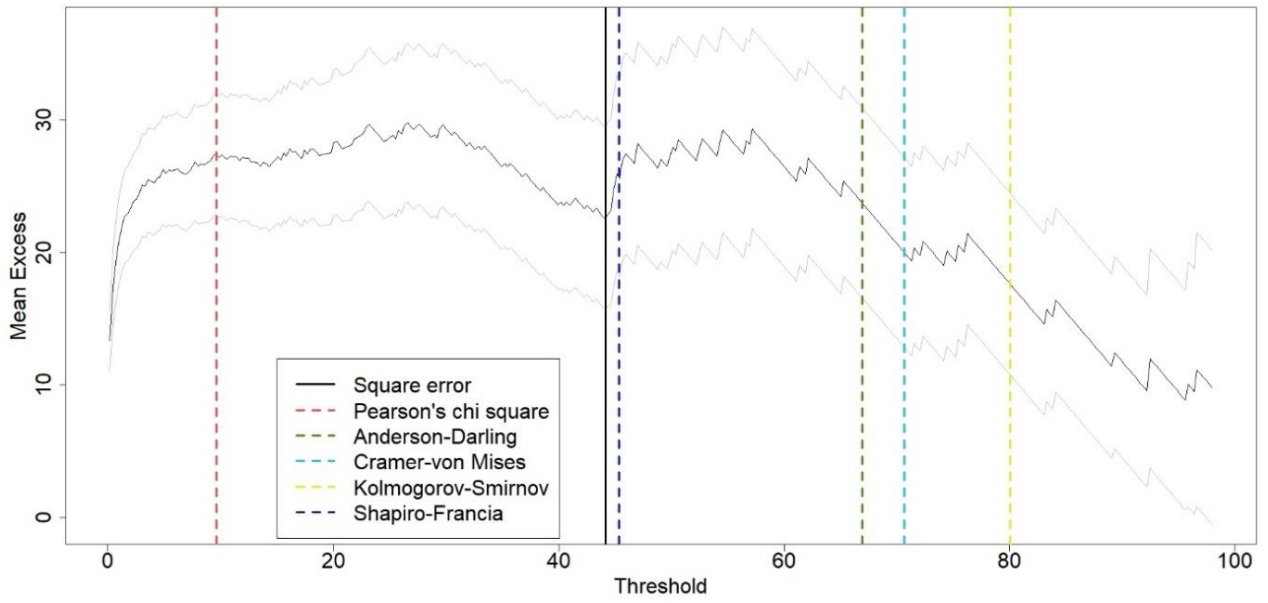
408 a)



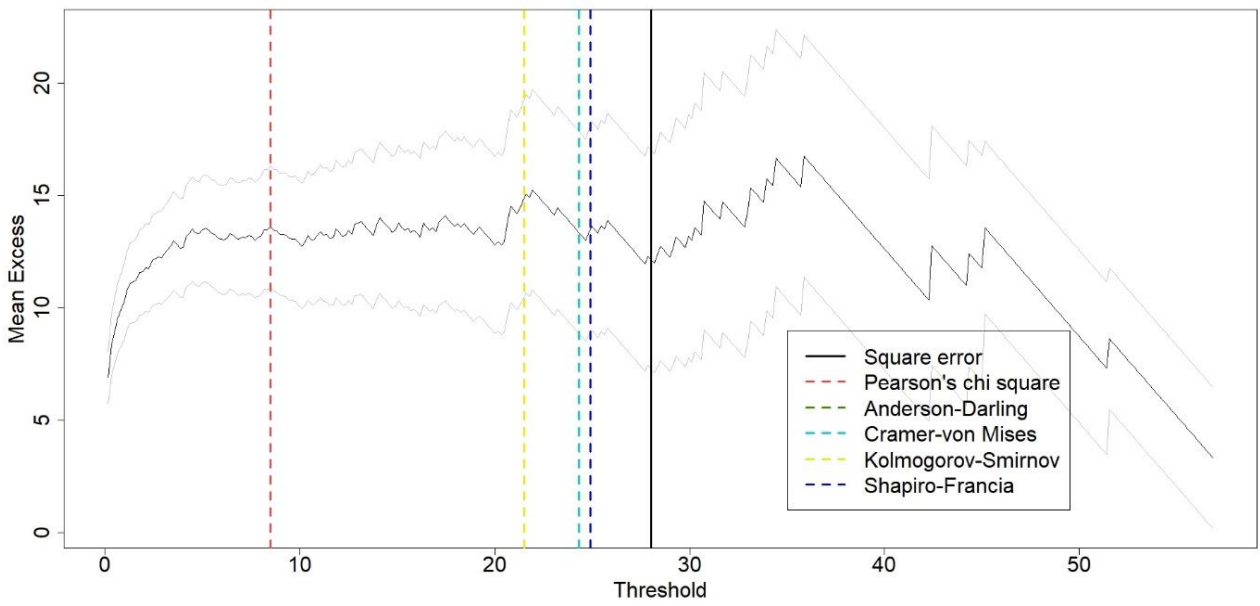
409

410 b)

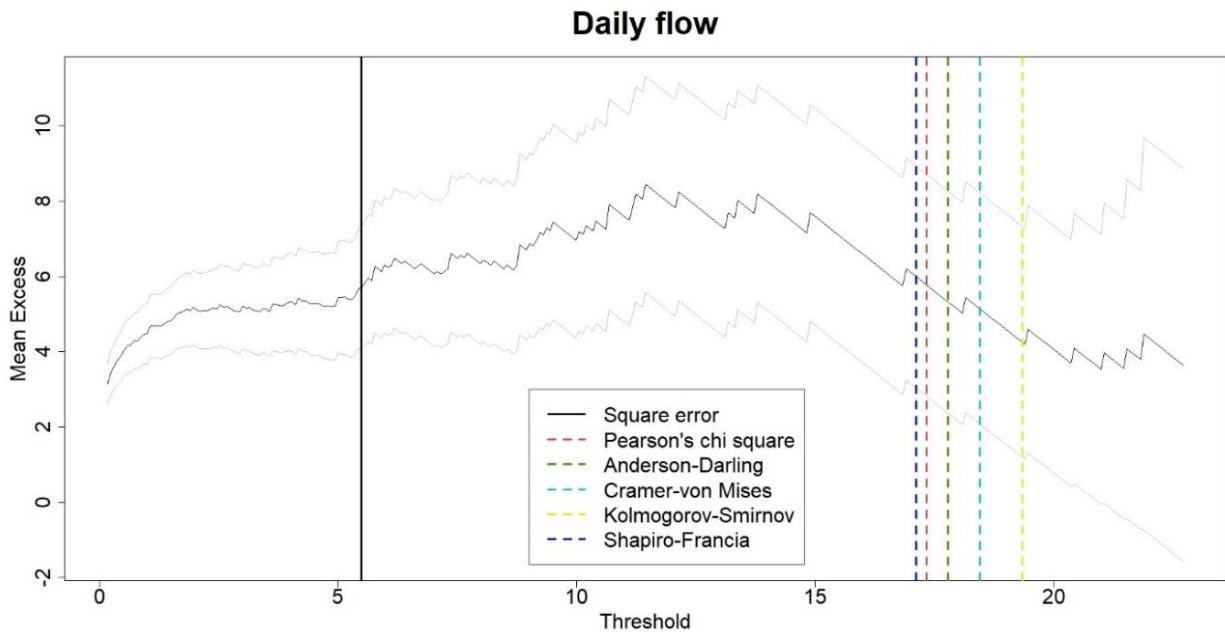
Hourly flow



6 hours flow



414 d)



415

416 Figure 7: MLR plots: Mean excesses and their 95% confidence intervals plotted against threshold for
417 the a) 15 minutes, b) hourly, c) 6 hourly and d) daily flow data. The threshold selected using the SE
418 method is shown by the vertical solid line and the thresholds selected by the Normality of
419 Differences tests are shown by the dashed vertical lines.

420 4.2.4 Parameter and fit comparisons

421 In summary, the estimated shape parameters showed little consistency across the four data
422 resolutions and across the threshold selection techniques investigated (Table 1). The 15-
423 minute extreme flows are characterized by: (i) an exponential tail (Pearson's chi square,
424 Anderson Darling and Kolmogorov-Smirnov tests) as the shape parameter takes values close
425 to zero, (ii) heavy tails (SE method, Shapiro-Francia and Cramer-von Mises tests) and (iii) short
426 tails ($\xi < 0$) (ATS method). ATS and Normality of Differences methods resulted in short tail
427 distributions for both the hourly and 6-hourly flow data, whereas the SE method resulted in
428 a heavier tail, similar to that found across all flow data scales. The ATS and the SE methods

429 provided heavy tails for the daily flow, and the Normality of Differences tests tended to short
 430 tails.

431 Table 1: Estimated thresholds and shape parameters for four flow resolutions and three core
 432 threshold selection methods.

		ATS	SE	Normality of Differences tests				
				Pearson's chi square	Anderson- Darling	Cramer- von Mises	Kolmogorov- Smirnov	Shapiro- Francia
15 mins	Threshold	22.2	46.8	39.7	51.8	45.5	42.6	53.5
	Shape Parameter	-0.14	0.33	0.01	0.07	0.26	0.06	0.10
Hourly	Threshold	9.7	44.7	9.6	66.9	70.7	80.1	45.3
	Shape Parameter	-0.09	0.17	-0.09	-0.58	-0.44	-0.48	-0.35
6 hours	Threshold	6.6	28.1	8.5	24.3	24.3	21.5	24.9
	Shape Parameter	-0.01	0.20	-0.05	-0.23	-0.23	-0.34	-0.23
Daily	Threshold	3.1	5.6	17.3	17.8	18.4	19.3	17.1
	Shape Parameter	0.17	0.22	-0.17	-0.10	-0.08	0.10	-0.20

433

434 Table 2: MSE between the empirical and theoretical quantiles for different threshold selection
 435 methods at four flow resolutions.

MSE	Threshold Stability	SE	Normality of Differences tests				
			Pearson's chi square	Anderson- Darling	Cramer- von Mises	Kolmogorov- Smirnov	Shapiro- Francia
15 mins	252.4	8248.8	123.7	2157.8	6034.9	1242.3	2828.2
Hourly	130.9	2654.1	24.1	14.5	13.6	10.5	28.0
6 hourly	72.1	150.8	61.0	34.0	34.0	12.7	34.8
Daily	38.2	81.9	8.3	10.7	12.6	32.4	7.6

436

437 The MSE (Table 2) seems to be an inappropriate diagnostic for deviations between very large
 438 theoretical and empirical quantiles as it depends greatly on the shape parameter. Peak flows
 439 with very short finite tails will show minimum MSEs, which increase by orders of magnitude

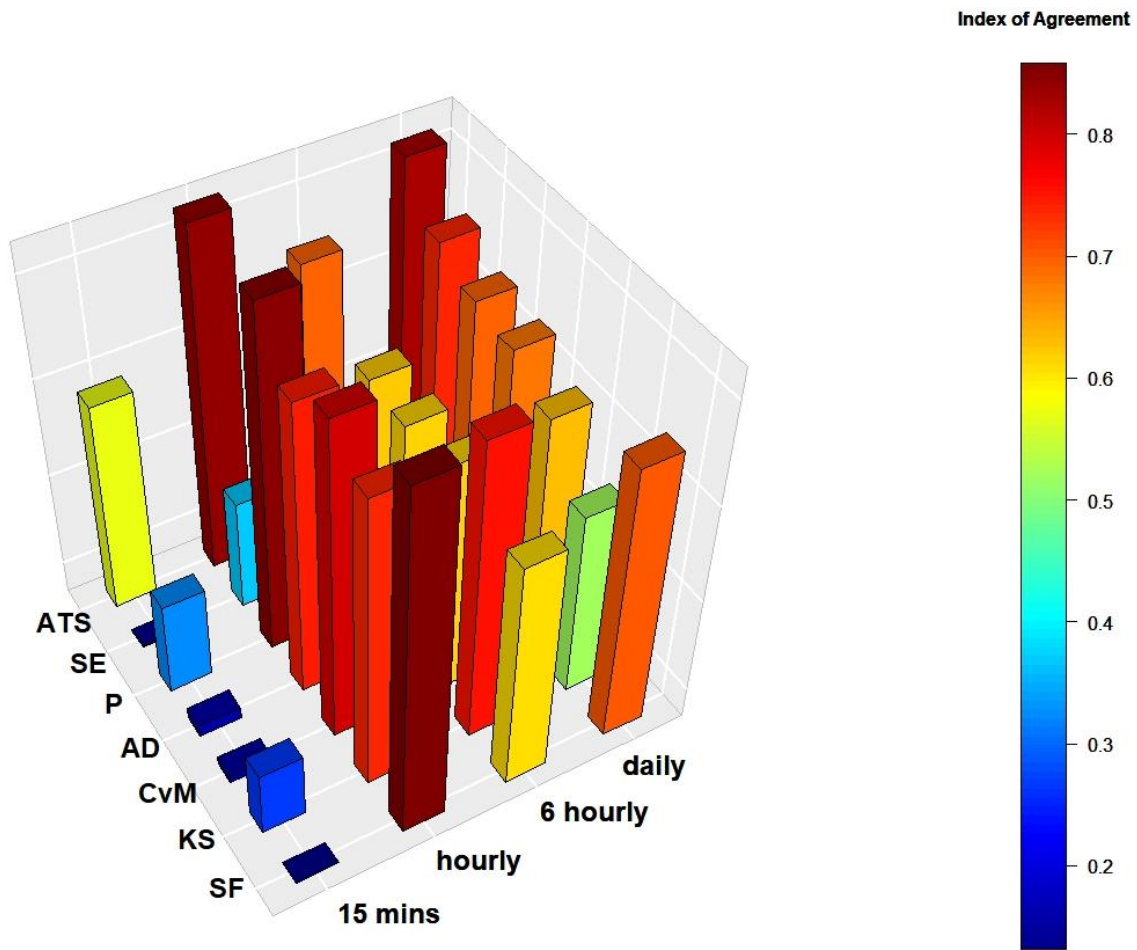
440 as the shape parameter increases. Conversely, the NRMSE does provide a comparative
 441 diagnostic since it is normalized by accounting for very large values that are associated with
 442 heavy tails. Thus, NRMSE values are reported in Table 3 where compared to the SE and
 443 Normality of Differences methods, this study's ATS method gives the smallest NRMSE for flow
 444 data of any resolution, except for the Normality of Differences test for the hourly flow.

445 Table 3: NRMSE between the empirical and theoretical quantiles for different threshold selection
 446 methods at four flow resolutions.

NRMSE	ATS	SE	Normality of Differences tests				
			Pearson's chi square	Anderson- Darling	Cramer- von Mises	Kolmogorov- Smirnov	Shapiro- Francia
15 mins	102.6	1017.9	308.0	571.6	866.6	391.4	697.5
Hourly	38.8	244.4	37.7	30.9	29.9	38.2	27.0
6 hourly	51.8	184.2	67.6	87.4	87.4	53.4	88.5
Daily	44.5	69.3	52.6	59.5	72.0	115.3	50.2

447

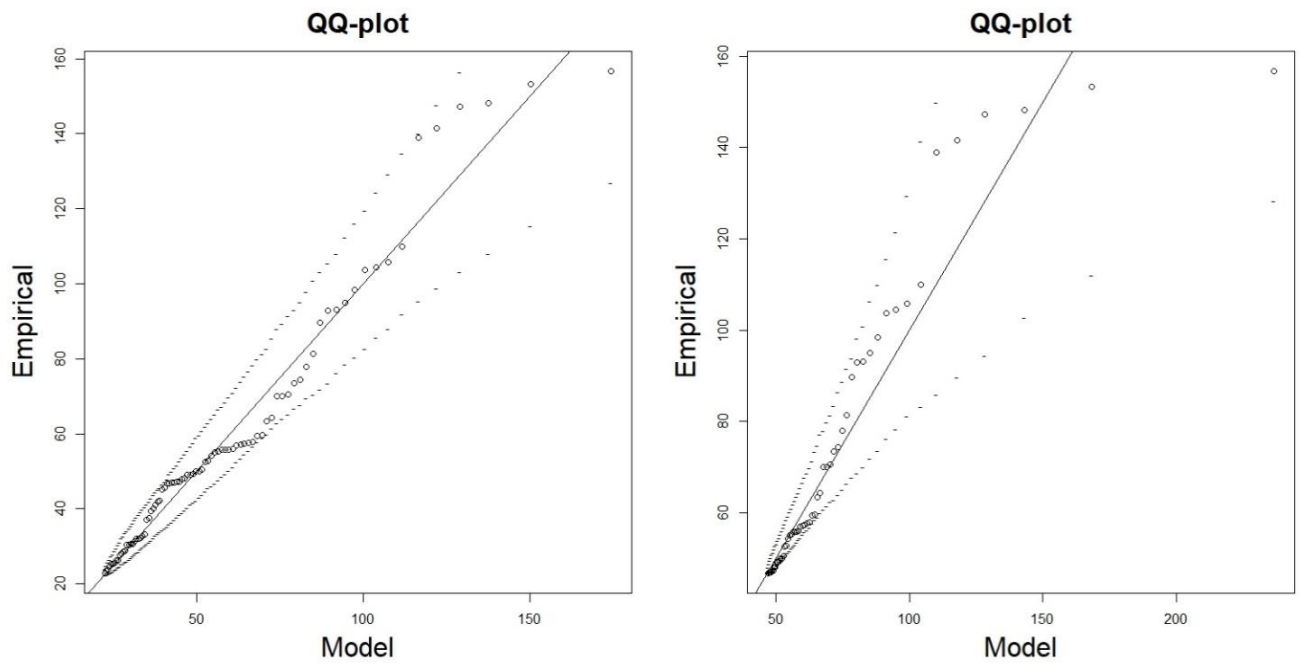
448 The relative index of agreement (Figure 8) is also an efficient measure of proximity between
 449 observed and simulated peak flows (Krause et al., 2005). For this diagnostic, the GPD was
 450 consistently best fitted to empirical peak flows at all scales when their thresholds were chosen
 451 using this study's ATS method. Here, the SE method was the poorest method, especially at
 452 the 15-minute data scale. Interestingly, results at the hourly scale behaved very differently to
 453 those found at the three other scales. We speculate that this was likely due to the hourly data
 454 being at, or close to, the natural water run-off integration rate to the sub-catchment's water
 455 flume following a rainfall event (see Discussion).



456

457 Figure 8: Index of agreement between theoretical and empirical peak flow of different resolutions.
 458 The threshold selection methods are Automated Threshold Stability (ATS), Square Error (SE) and the
 459 various tests of the Normality of Differences method, the Pearson's chi-square (P), Anderson-Darling
 460 (AD), Cramer-von Mises (CvM), Kolmogorov-Smirnov (KS) and Shapiro-Francia (SF).

461 Figure 9 presents the Q-Q plots of the 15-minute extreme flows for the threshold selection
 462 methods that gave the smallest (ATS) and the largest (SE) NRMSE values (Table 3). The Q-Q
 463 plots show that an over-estimated threshold results in a sample size that can be too small for
 464 efficient statistical inference and results in increased uncertainty. The Q-Q plots also emphasis
 465 the superiority of this study's ATS method given its Q-Q plot falls relatively close to the 45°
 466 line.

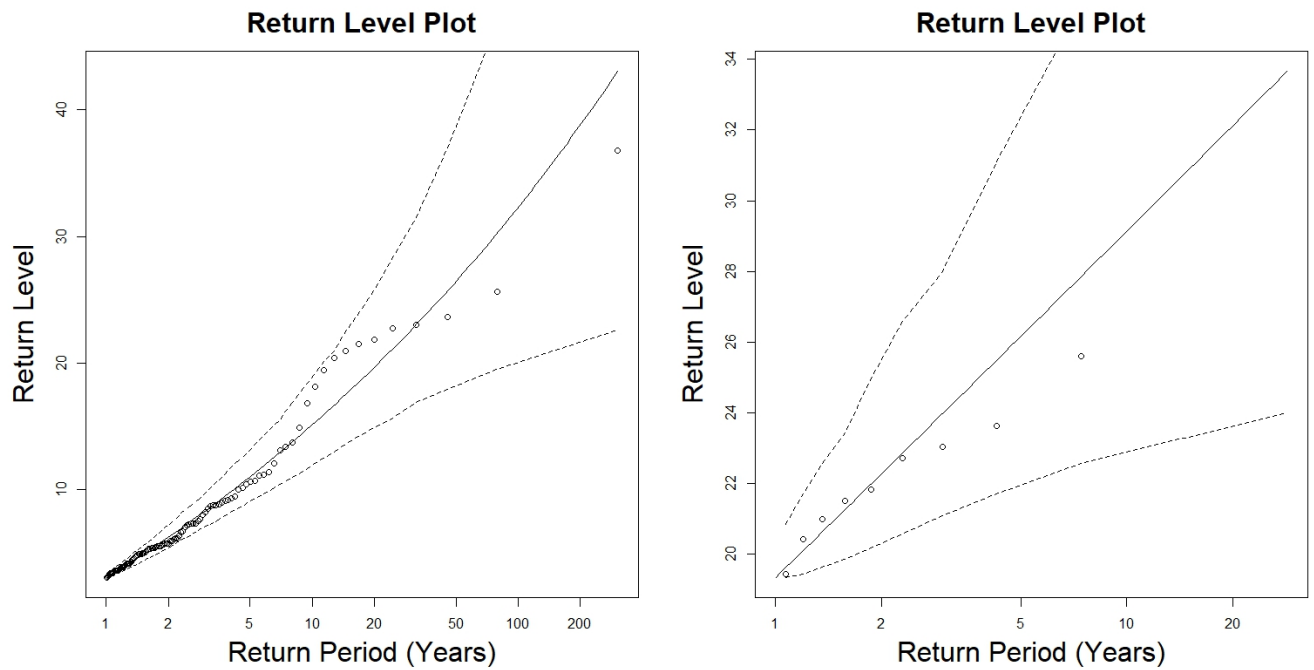


467

468 Figure 9: Q-Q plots of the 15-minute peak flows estimated by the ATS (left) and SE (right) methods.

469

470 Clear differences in the estimated Return Level / Return Period plots for the ATS and
 471 Normality of Difference (Kolmogorov-Smirnov test only) methods (Figure 10) indicate that the
 472 combined effects of data scale, the GPD estimator and the threshold selection method - each
 473 have a significant impact on the characteristics of the final model that attempts to explain the
 474 flow process with the consideration of extremes. This is critically important in cases where
 475 reliably informed actions need to be taken or infrastructure needs to be built to mitigate the
 476 impacts of future peak flows and likely flood events.



477

478 Figure 10: Return level plots of the daily peak flows estimated by the ATS (left) and Normality of
 479 Difference Kolmogorov-Smirnov (right) methods.

480 5. Discussion

481 In agreement with previous studies (e.g. Bermudez & Kotz, 2010; Engeland et al., 2004), we
 482 found that the performance of the GPD parameter estimators examined through a Monte
 483 Carlo experiment, depended significantly on the sample size and the value of the shape
 484 parameter. The MLE/MPLE, PWMU/PWMB and the LME were consistently the most unbiased
 485 and precise estimators and so we chose only from this group in our subsequent analyses.
 486 More specifically, for the application of the SE and AST threshold selection methods, a
 487 different GPD estimator was used each time according to its strengths. For example, the LME
 488 was preferred for positive shape parameters and large sample size.

489 This study's Automated Threshold Stability (ATS) method was tested against existing SE and
 490 Normality of Differences methods. Methods were applied to flow discharge measurements
 491 of 15-minute resolution, as well as to the same data aggregated to coarser resolutions of

492 hourly, 6-hourly and daily, to examine scale effects. The Normality of Differences method
493 depended on the normality test applied and resulted in short, exponential and heavy tailed
494 distributions even at the same scale (e.g. shape parameters of $\xi = -0.2$ for the daily flow
495 according to Shapiro-Francia and $\xi = 0.1$ according to the Kolmogorov-Smirnov test). Similar
496 results for the value of the shape parameter were obtained from the ATS method, unlike the
497 SE method which always resulted in positive ξ .

498 Threshold stability plots were discussed in Scarrott and MacDonald (2012) and Solari and
499 Losada (2012), but these studies did not perform an analytical approximation, as done here
500 with ATS, although Langousis et al. (2016) suggested an automated technique based on the
501 assumption of linearity of the MRL plot and applied it to rainfall data. Our proposed ATS
502 method provided more robust estimates of the threshold compared to: (a) the SE method as
503 it was less sensitive to the resolution of the data and (b) the Normality of Differences method
504 as it was less sensitive to the sample size of the threshold candidates. It also resulted in the
505 smallest errors and the largest agreement indices between the simulated and the empirical
506 quantiles.

507 Specific to the case study, error and agreement indices indicated that the GPD provided the
508 best fit to the hourly peak flow data relative to 15-minute, 6-hourly and daily peak flow data.
509 For all the applied threshold selection methods, the modelled peak flow at the hourly
510 resolution was consistently the closest to the empirical one, compared to three other scales.
511 These results cannot be attributed to the value of the shape parameter (e.g. short finite tails
512 result in greater agreement between theoretical and empirical quantiles) since the SE method
513 gives a positive ξ . An inspection of the plots and a comparison across various scales does not
514 reveal any pattern that would justify this behavior. A possible explanation could be that the

515 hourly peak flow best captures the signal of the process and integrates more efficiently the
516 way the 6.84 ha sub-catchment (of two pasture fields) transforms intensive rainfall into high
517 discharge flows. It should be noted that the data aggregation was not done at equal intervals.
518 For example, the hourly flow resulted from averaging four 15-minute measurements,
519 whereas the 6-hourly and the daily flow are the averages of 24 and 96 observations,
520 respectively. This does not affect the results but should be borne in mind when interpreting
521 the plots.

522 An advantage of using fine resolution flow data is that they result in larger sample sizes that
523 can make the statistical inference more efficient even for records of short periods for which
524 a GEV/AM extreme value methodology is not applicable. However, this study showed that for
525 data of the same resolution, the value of the GDP shape parameter varies according to the
526 selected thresholds. This has serious practical implications since the models are commonly
527 extrapolated beyond observed values for forecasting and engineering design purposes to
528 mitigate against future flooding. On one hand, an under-estimated threshold and shape
529 parameter of the extreme flow can result in failure of hydrological infrastructure (e.g. dams,
530 flood protection works) due to higher peak flows than expected. On the other hand, over-
531 estimation of the high flows can lead to over-pricing and mis-use of resources.

532 6. Conclusions

533 In this study, we examined the effect of statistical estimators, data resolution, and threshold
534 selection on fitting the Generalized Pareto distribution to peak hydrological flows that
535 resulted from the 'Peaks Over Threshold' method. Through a simulation study, the
536 performance of the estimators depended greatly on the sample size and the shape parameter

537 where the only most accurate and unbiased estimators were used for the selection of
538 thresholds in subsequent empirical evaluations. Here an automated threshold selection
539 method based on the stability of the shape and modified scale parameters was empirically
540 demonstrated to provide more robust estimates compared to two commonly applied
541 alternatives. The proposed method provided the smallest error and the greatest agreement
542 indices between the empirical and theoretical quantiles across all the scales of the case study
543 flow data.

544 The study results can be generalized to similar water monitoring schemes for improved
545 characterization of likely flood events. However, the study highlights that the combined effect
546 of data scale, threshold selection method and statistical estimator, significantly affects the
547 shape parameter and, as a consequence, the nature of the Generalized Pareto distribution.
548 Such linked effects need to be acknowledged and assessed as they have clear implications for
549 the reliable forecasting of extreme flow events, and the consequences thereof.

550 ***Acknowledgements***

551 Rothamsted Research receives grant aided support from the Biotechnology and Biological
552 Sciences Research Council (BBSRC) of the United Kingdom. This research was funded by
553 Rothamsted Research and Lancaster Environment Centre, the BBSRC Institute Strategic
554 Programme (ISP) grant, “Soils to Nutrition” (S2N) grant numbers BBS/E/C/00010320,
555 BBS/E/C/00010330 and the BBSRC National Capability grant for the North Wyke Farm Platform
556 grant number BBS/E/C/000J0100.

557 ***Declaration of interest***

558 The authors declare no potential conflict of interest associated with this research.

559 ***Software and data availability***

560 The statistical software (R Core Team, 2017) and all North Wyke Farm Platform data sets
561 (<https://www.rothamsted.ac.uk/north-wyke-farm-platform>) are freely available.

562

563 References

- 564 Ashkar, F. and Tatsambon, C. N. (2007). Revisiting some estimation methods for the
565 generalized Pareto distribution. *Journal of Hydrology*, 346, 136-143.
- 566 Bates, B. C., Kundzewicz, Z. W., Wu, S. and Palutikof, J. P. (2008). *Climate Change and*
567 *Water*. Technical Paper of the Intergovernmental Panel on Climate Change, IPCC
568 Secretariat, Geneva, 210 pp.
- 569 Beguería, S. (2005). Uncertainties in Partial Duration Series Modelling of Extremes Related to
570 the Choice of the Threshold Value. *Journal of Hydrology*, 303(1), 215-230.
- 571 Behrens, C. N., Lopes, H. F. and Gamerman, D. (2004). Bayesian Analysis of Extreme Events
572 with Threshold Estimation. *Statistical Modelling*, 4(3), 227-244.
- 573 Beirlant, J., Dierckx, G. and Guillou, A. (2005). Estimation of the Extreme-Value Index and
574 Generalized Quantile Plots. *Bernoulli*, 11(6), 949-970.
- 575 Beirlant, J., Joossens, E. and Segers, J. (2005). Unbiased tail estimation by an extension of the
576 generalized Pareto distribution. *CentER Discussion Paper*, Vol. 2005-112, Tilburg:
577 *Econometrics*.
- 578 Beirlant, J., Vynckier, P. and Teugels, J. L. (1996). Tail Index Estimation, Pareto Quantile
579 Plots, and Regression Diagnostics. *Journal of the American Statistical Association*,
580 91(436), 1659-1667.
- 581 Beirlant, J., de Wet, T. and Goegebeur, Y. (2006). A Goodness-of-Fit Statistic for Pareto-Type
582 Behaviour. *Journal of Computational and Applied Mathematics*, Special Issue: Jef
583 Teugels, 186(1), 99-116.
- 584 Bommier, E. (2014). *Peaks-over-threshold modelling of environmental data (Technical report)*.
585 U.U.D.M. Project Report, 2014:33.

586 Bouraoui, F., Grizzetti, B., Granlund, K., Rekolainen, S. and Bidoglio, G. (2004). Impact of
587 Climate Change on the Water Cycle and Nutrient Losses in a Finnish Catchment,
588 Climatic Change, 66(1–2), 109-126.

589 Brodin, E. and Rootzén, H. (2009). Univariate and Bivariate GPD Methods for Predicting
590 Extreme Wind Storm Losses. Insurance: Mathematics and Economics, 44(3), 345-356.

591 Choulakian, V. and Stephens, M. A. (2001). Goodness-of-Fit Tests for the Generalized Pareto
592 Distribution. Technometrics, 43(4), 478-484.

593 Claps, P. and F. Laio. (2003). Can Continuous Streamflow Data Support Flood Frequency
594 Analysis? An Alternative to the Partial Duration Series Approach. Water Resources
595 Research, 39(8).

596 Clarke, M. L. & Rendell, H. M. (2006). Hindcasting Extreme Events: The Occurrence and
597 Expression of Damaging Floods and Landslides in Southern Italy. Land Degradation &
598 Development, 17(4), 365-380.

599 Coles, S. (2001). An Introduction to Statistical Modeling of Extreme Values. Springer, London,
600 UK.

601 Coles, S. and Dixon, M. J. (1999). Likelihood-Based Inference for Extreme Value Models,
602 Extremes, 2(1), 5-23.

603 Cunnane, C. (1979). A Note on the Poisson Assumption in Partial Duration Series Models.
604 Water Resources Research, 15(2), 489-494.

605 Danielsson, J., de Haan, L., Peng, L. and de Vries, C. G. (2001). Using a Bootstrap Method to
606 Choose the Sample Fraction in Tail Index Estimation. Journal of Multivariate Analysis,
607 76(2), 226-248.

608 Das, B. and Ghosh, S. (2013). Weak limits for exploratory plots in the analysis of extremes.
609 Bernoulli, 19(1), 308-343

610 Davison, A. C. and Smith, R. L. (1990). Models for Exceedances over High Thresholds,
611 Journal of the Royal Statistical Society, Series B (Methodological), 52(3), 393-442.

612 Deidda, R. (2010). A Multiple Threshold Method for Fitting the Generalized Pareto
613 Distribution to Rainfall Time Series. Hydrology and Earth System Sciences, 14(12),
614 2559-2575.

615 Deidda, R. and Puliga, M. (2006). Sensitivity of Goodness-of-Fit Statistics to Rainfall Data
616 Rounding Off. Physics and Chemistry of the Earth, 31(18). 1240-1251.

617 Dekkers, A. L. M. and De Haan, L. (1989). On the Estimation of the Extreme-Value Index and
618 Large Quantile Estimation, The Annals of Statistics, 17(4), 1795-1832.

619 Durocher, M., Zadeh, S. M., Burn, D. H. and Ashkar, F. (2018). Comparison of Automatic
620 Procedures for Selecting Flood Peaks over Threshold Based on Goodness-of-Fit Tests.
621 Hydrological Processes, 32(18), 2874–2887.

622 Eastoe, E. F. and Tawn, J. A. (2010). Statistical Models for Overdispersion in the Frequency
623 of Peaks over Threshold Data for a Flow Series. Water Resources Research, 46(2).

624 Engeland, K., Hisdal, H. and Frigessi, A. (2004). Practical Extreme Value Modelling of
625 Hydrological Floods and Droughts: A Case Study. Extremes, 7(1), 5–30.

626 Ferguson, T. S., Genest, C. and Hallin, M. (2000). Kendall’s Tau for Serial Dependence.
627 Canadian Journal of Statistics, 28(3), 587–604.

628 Field, C. B., Barros, V., Stocker, T. F. and Dahe, Q. (2012). Managing the Risks of Extreme
629 Events and Disasters to Advance Climate Change Adaptation: Special Report of the
630 Intergovernmental Panel on Climate Change, Cambridge, Cambridge University Press.

631 Fisher, R. A. and Tippett, L. H. C. (1928). Limiting forms of the frequency distribution of the
632 largest or smallest member of a sample, *Proc. Cambridge Philos. Soc.*, 24(2), 180-190.

633 Gharib, A., Davies, E. G. R., Goss, G. G. and Faramarzi, M. (2017). Assessment of the
634 Combined Effects of Threshold Selection and Parameter Estimation of Generalized
635 Pareto Distribution with Applications to Flood Frequency Analysis, *Water*, 9(9), 692.

636 Goegebeur, Y., Beirlant, J. and de Wet, T. (2008). Linking Pareto-Tail Kernel Goodness-of-
637 Fit Statistics with Tail Index at Optimal Threshold and Second Order Estimation.,
638 *REVSTAT-Statistical Journal*, 6(1), 51-69.

639 Greenwood, J. A., Landwehr, J. M., Matalas N. C. and Wallis, J. R. (1979). Probability
640 Weighted Moments: Definition and Relation to Parameters of Several Distributions
641 Expressable in Inverse Form. *Water Resources Research*, 15(5), 1049-1054.

642 Hall, P. (1990). Using the Bootstrap to Estimate Mean Squared Error and Select Smoothing
643 Parameter in Nonparametric Problems. *Journal of Multivariate Analysis*, 32(2), 177-
644 203.

645 Hill, B. M. (1975). A Simple General Approach to Inference About the Tail of a Distribution.
646 *The Annals of Statistics*, 3(5), 1163-1174.

647 Hosking, J. R. M. and Wallis, J. R. (1987). Parameter and Quantile Estimation for the
648 Generalized Pareto Distribution. *Technometrics*, 29(3), 339-349.

649 Jenkinson, A. F. (1955). The Frequency Distribution of the Annual Maximum (or Minimum)
650 Values of Meteorological Elements. *Quarterly Journal of the Royal Meteorological*
651 *Society*, 81(348), 158-171.

652 Josse, J. and Husson, F. (2013). Handling Missing Values in Exploratory Multivariate Data
653 Analysis Methods. *Journal de La Société Française de Statistique*, 153(2), 79–99.

654 Krause, P., Boyle, D. P. and Bäse, F. (2005). Comparison of Different Efficiency Criteria for
655 Hydrological Model Assessment. *Advances in Geosciences*, 5, 89–97.

656 Kundzewicz, Z. W., Mata, L. J., Arnell, N. W., Doll, P., Kabat, P., Jimenez, B. et al. (2007).
657 Freshwater Resources and Their Management. In *Climate Change 2007: Impacts,*
658 *Adaptation and Vulnerability. Contribution of Working Group II to the Fourth*
659 *Assessment Report of the Intergovernmental Panel on Climate Change*, edited by M.
660 L. Parry, O. F. Canziani, J. P. Palutikof, P. J. van der Linden, and C. E. Hanson, 173–
661 210. Cambridge University Press.

662 Landwehr, J. M., Matalas, N. C. and Wallis, J. R. (1979). Probability Weighted Moments
663 Compared with Some Traditional Techniques in Estimating Gumbel Parameters and
664 Quantiles. *Water Resources Research*, 15(5), 1055-1064.

665 Lang, M., Ouarda, T. B. M. J. and Bobée, B. (1999). Towards Operational Guidelines for Over-
666 Threshold Modeling. *Journal of Hydrology*, 225(3), 103-117.

667 Langousis, A., Mamalakis, A., Puliga, M. & Deidda, R. (2016). Threshold Detection for the
668 Generalized Pareto Distribution: Review of Representative Methods and Application
669 to the NOAA NCDC Daily Rainfall Database. *Water Resources Research*, 52(4), 2659–
670 2681.

671 Ledford, A. W. and Tawn, J. A. (2003). Diagnostics for Dependence within Time Series
672 Extremes. *Journal of the Royal Statistical Society. Series B (Statistical Methodology)*,
673 65(2), 521–543.

674 Liang, B, Shao, Z., Li, H., Shao, M. and Lee, D. (2019). An Automated Threshold Selection
675 Method Based on the Characteristic of Extrapolated Significant Wave Heights. *Coastal*
676 *Engineering*, 144, 22-32.

677 Liu, Y., Li, Y., Harris, P., Cardenas, L. M., Dunn, R. M., Sint, H., Murray, P. J., Lee, M. R. F.
678 and Wu, L. (2018). Modelling Field Scale Spatial Variation in Water Run-off, Soil
679 Moisture, N₂O Emissions and Herbage Biomass of a Grazed Pasture Using the
680 SPACSYS Model. *Geoderma*, 315, 49-58.

681 Luceño, A. (2006). Fitting the Generalized Pareto Distribution to Data Using Maximum
682 Goodness-of-Fit Estimators. *Computational Statistics & Data Analysis*, 51(2), 904-917.

683 Mackay, E. B. L., Challenor, P. G. and Bahaj, A. S. (2011). A Comparison of Estimators for
684 the Generalised Pareto Distribution. *Ocean Engineering*, 38(11), 1338-1346.

685 Madsen, H., Rasmussen, P. F. and Rosbjerg, D. (1997). Comparison of Annual Maximum
686 Series and Partial Duration Series Methods for Modeling Extreme Hydrologic Events:
687 1. At-Site Modeling. *Water Resources Research*, 33(4), 747–757.

688 Millennium Ecosystem Assessment, 2005. Millennium Ecosystem Assessment (MA), Ecosystems
689 and Human Well-being: Synthesis. Island Press, Washington, DC.

690 Moore, D. S. (1986). Tests of Chi-Squared Type Goodness of Fit Techniques. Marcel Dekker,
691 New York.

692 Orr, R. J., Murray, P. J., Eyles, C. J., Blackwell, M. S. A., Cardenas, L. M., Collins, A. L. et al.
693 (2016). The North Wyke Farm Platform: effect of temperate grassland farming systems
694 on soil moisture contents, runoff and associated water quality dynamics, *European*
695 *Journal of Soil Science*, 67, 374–385.

696 de Zea Bermudez, P. and Kotz, S. (2010). Parameter Estimation of the Generalized Pareto
697 Distribution, Part I. *Journal of Statistical Planning and Inference*, 140(6), 1353-1373.

698 Pickands, J. (1975). Statistical Inference Using Extreme Order Statistics. *The Annals of*
699 *Statistics*, 3(1), 119-131.

700 Prescott, P. and Walden, A. T. (1980). Maximum Likelihood Estimation of the Parameters of
701 the Generalized Extreme-Value Distribution. *Biometrika*, 67(3), 723-724.

702 Prescott, P. and Walden, A. T. (1983). Maximum Likelihood Estimation of the Parameters of
703 the Three-Parameter Generalized Extreme-Value Distribution from Censored Samples.
704 *Journal of Statistical Computation and Simulation*, 16(3–4), 241-250.

705 R Core Team (2017). R: A language and environment for statistical computing. R Foundation
706 for Statistical Computing, Vienna, Austria. URL <https://www.R-project.org/>.

707 Reiss, R. D. and Thomas, M. (2007). *Statistical Analysis of Extreme Values: With Applications*
708 *to Insurance, Finance, Hydrology and Other Fields*. 3rd edition. Birkhäuser Basel.

709 Scarrott, C. and MacDonald, A. (2012). A Review of Extreme Value Threshold Es-Timation
710 and Uncertainty Quantification. *REVSTAT–Statistical Journal*, 10(1), 33–60.

711 Segers, J. (2005). Generalized Pickands Estimators for the Extreme Value Index. *Journal of*
712 *Statistical Planning and Inference*, 128(2), 381-396.

713 Sheta, A. F. and El-Sherif, M. S. (1999). Optimal Prediction of the Nile River Flow Using
714 Neural Networks. *International Joint Conference on Neural Networks. Proceedings*, 5,
715 3438-3441.

716 Sigauke, C. and Bere, A. (2017). Modelling Non-Stationary Time Series Using a Peaks over
717 Threshold Distribution with Time Varying Covariates and Threshold: An Application
718 to Peak Electricity Demand. *Energy*, 119, 152-166.

719 Smith, R. L. (1985). Maximum Likelihood Estimation in a Class of Nonregular Cases.
720 *Biometrika*, 72(1), 67-90.

721 Solari, S. and Losada, M. A. (2012). A Unified Statistical Model for Hydrological Variables
722 Including the Selection of Threshold for the Peak over Threshold Method. *Water*
723 *Resources Research*, 48(10).

724 Solari, S., Egüen, M., Polo, M. J. and Losada, M. A. (2017). Peaks Over Threshold (POT): A
725 Methodology for Automatic Threshold Estimation Using Goodness of Fit p-Value.
726 *Water Resources Research*, 53(4), 2833–2849.

727 Takahashi, T., Harris, P., Blackwell, M. S. A., Cardenas, L. M., Collins, A. L., Dungait, J. A.
728 J. et al. (2018). Roles of Instrumented Farm-Scale Trials in Trade-off Assessments of
729 Pasture-Based Ruminant Production Systems. *Animal*, 12 (8), 1766-1776.

730 Tanaka, S. and Takara, K. (2002). A Study on Threshold Selection in POT Analysis of Extreme
731 Floods, Extremes of the Extremes, Extraordinary Floods, *IAHS Publ*, 271, 299-304.

732 Thibault, K. M. & Brown, J. H. (2008). Impact of an Extreme Climatic Event on Community
733 Assembly. *Proceedings of the National Academy of Sciences*, 105(9), 3410-3415.

734 Thode, H. C. (2002). *Testing For Normality*. CRC Press.

735 Thompson, P., Cai, Y., Reeve, D. and Stander, J. (2009). Automated Threshold Selection
736 Methods for Extreme Wave Analysis. *Coastal Engineering*, 56(10), 1013-1021.

737 Todorovic, P. (1978). Stochastic Models of Floods. *Water Resources Research*, 14(2), 345–
738 356.

739 Turan, M. E., and Yurdusev, M. A. (2009). River Flow Estimation from Upstream Flow
740 Records by Artificial Intelligence Methods. *Journal of Hydrology*, 369(1), 71-77.

741 Willmott, C. J. (1981). On the Validation of Models. *Physical Geography*, 2(2), 184-194.

742 Yang, X., Zhang, J. and Ren, W. X. (2018). Threshold Selection for Extreme Value Estimation
743 of Vehicle Load Effect on Bridges. *International Journal of Distributed Sensor*
744 *Networks*, 14(2).

745 Yun, S. (2002). On a Generalized Pickands Estimator of the Extreme Value Index. *Journal of*
746 *Statistical Planning and Inference*, 102(2), 389-409.

747 de Zea Bermudez, P. and Kotz, S. (2010). Parameter Estimation of the Generalized Pareto
748 Distribution, Part I. *Journal of Statistical Planning and Inference*, 140(6), 1353–1373.

749 Zhang, J. (2007). Likelihood Moment Estimation for the Generalized Pareto Distribution.
750 *Australian & New Zealand Journal of Statistics*, 49(1), 69-77.

751 Zoglat, A., EL Adlouni, S., Badaoui, F., Amar A. & Okou, C. G. (2014). Managing
752 Hydrological Risks with Extreme Modeling: Application of Peaks over Threshold
753 Model to the Loukkos Watershed, Morocco, *Journal of Hydrologic Engineering*, 19(9),
754 05014010.

755

756

757 Appendix A: Equations of the estimators

758 The estimators used in this study can be formally defined as follows:

759 1. MLE method:

$$760 \quad L = -n \log \sigma + \left(\frac{1}{\xi} - 1 \right) \sum_{i=1}^n \log \left(1 - \frac{\xi x_i}{\sigma} \right), \quad \xi \neq 0$$

$$761 \quad L = -n \log \sigma - \frac{1}{\sigma} \sum_{i=1}^n x_i, \quad \xi = 0$$

762 where $x_{(1)} \leq x_{(2)} \leq \dots \leq x_{(n)}$ are the order statistics of a random sample x_1, \dots, x_n from the
 763 GPD. The estimated parameters are obtained when the log-likelihood function L is
 764 maximized.

765 2. MPLE method:

$$766 \quad P(\xi) = \begin{cases} 1 & \xi \leq 0 \\ \exp\{-\lambda \left(\frac{1}{1-\xi} - 1 \right)^a\} & 0 < \xi < 1 \\ 0 & \xi \geq 1 \end{cases}$$

767 where a and λ are the penalizing non-negative constants. The corresponding penalized
 768 likelihood function is $L_{pen} = L \times P$.

769 3. LME is a combination of both likelihood and moment estimators and is derived from:

$$770 \quad \frac{1}{n} \sum_{i=1}^n (1 - \theta x_i)^P - \frac{1}{1-r} = 0, \quad \theta < x_{(n)}^{-1},$$

771 where $\theta = \xi/\sigma$ and $P = -\frac{rn}{\sum_{i=1}^n \log(1-\theta x_i)}$. The parameter $r < 1, r \neq 0$ must be pre-defined

772 before the estimation and either be set as ξ if there is an initial estimate of it or taken as

773 $r = -1/2$.

774 4. MOM estimators (Hosking & Wallis, 1987) of the scale σ and shape ξ parameters of the
 775 GPD distribution are given by:

$$776 \quad \hat{\sigma} = \frac{1}{2}\bar{x}\left(\frac{\bar{x}}{s^2} + 1\right), \quad \hat{\xi} = \frac{1}{2}\left(\frac{\bar{x}^2}{s^2} - 1\right)$$

777 where \bar{x} and s^2 are the sample mean and variance.

778 5. PWM estimators provide estimates with smaller bias and variance than MLE when the
 779 sample size is less than 500 (Hosking & Wallis 1987). The PWM's of the random variable
 780 X with a distribution function $G \equiv G(x) = P(X \leq x)$ is defined as:

$$781 \quad M_{l,j,k} = E[X^l F^j (1 - F)^k] = \int_0^1 [x(F)]^l F^j (1 - F)^k dF$$

782 where l, j and k are real numbers. For $j = k = 0$ and l a nonnegative integer, $M_{l,0,0}$ is the
 783 classical moment of order l .

784 6. The estimator suggested by Pickands (1975) (referred to as 'Pick') is based on the
 785 ascending order statistics $X_{1,n} \leq X_{2,n} \leq \dots \leq X_{n,n}$ from an independent sample of size n
 786 and is defined as:

$$787 \quad \hat{\xi}_{n,k}^{Pick} = \frac{1}{\log 2} \log \left(\frac{X_{n-k+1,n} - X_{n-2k+1,n}}{X_{n-2k+1,n} - X_{n-4k+1,n}} \right), \text{ for } k = 1, \dots, [n/4]$$

788 This estimator is largely dependent on k and provides a large asymptotic variance (e.g.
 789 (Dekkers & Haan, 1989; Segers, 2005; Yun, 2002).

790 7. There are many MGF statistics that can be used for GPD parameter estimation, such as
 791 Kolmogorov-Smirnov, Cramer-von Mises and Anderson-Darling (see Luceño, 2006).

792

793
794
795
796
797
798
799
800

List of Tables

Table 1: Estimated thresholds and shape parameters for four flow resolutions and three core threshold selection methods.	26
Table 2: MSE between the empirical and theoretical quantiles for different threshold selection methods at four flow resolutions.	26
Table 3: NRMSE between the empirical and theoretical quantiles for different threshold selection methods at four flow resolutions.....	27

801
802
803
804
805
806
807
808
809
810
811
812
813
814
815
816
817
818
819
820
821
822

List of Figures

Figure 1: The three farmlets and the 15 sub-catchments of the North Wyke Farm Platform, with: (i) ‘blue’ farmlet a mixture of white clover and high sugar perennial ryegrass; (ii) ‘red’ farmlet high sugar perennial ryegrass only and (iii) ‘green’ farmlet permanent pasture (“business as usual”).13

Figure 2: Flow ($l\ s^{-1}$) measurements at sub-catchment 3 (2012 to 2018).14

Figure 3: Performance of GPD estimators for shape parameter $\xi = 0$ and for six different sample sizes ($n = 25, 50, 100, 250, 500, 1000$).16

Figure 4: Shape parameter characteristics of measured (15-minute) and a series of averaged (30-minute to daily) flow rates.18

Figure 5: Kendall’s test statistics τ (solid lines) along with the 95% acceptance limits of the test (dashed lines).18

Figure 6: Automated Threshold Stability (ATS) method: Selected threshold (that between the vertical green lines) of a) 15 minutes, b) hourly, c) 6 hourly and d) daily flow based on smoothing splines.21

Figure 7: MLR plots: Mean excesses and their 95% confidence intervals plotted against threshold for the a) 15 minutes, b) hourly, c) 6 hourly and d) daily flow data. The threshold selected using the SE method is shown by the vertical solid line and the thresholds selected by the Normality of Differences tests are shown by the dashed vertical lines.25

Figure 8: Index of agreement between theoretical and empirical peak flow of different resolutions. The threshold selection methods are Automated Threshold Stability (ATS), Square Error (SE) and the various tests of the Normality of Differences method, the Pearson’s

823	chi-square (P), Anderson-Darling (AD), Cramer-von Mises (CvM), Kolmogorov-Smirnov (KS)	
824	and Shapiro-Francia (SF).	28
825	Figure 9: Q-Q plots of the 15-minute peak flows estimated by the ATS (left) and SE (right)	
826	methods.	29
827	Figure 10: Return level plots of the daily peak flows estimated by the ATS (left) and Normality	
828	of Difference Kolmogorov-Smirnov (right) methods.	30
829		