

1 **AUTOMATED THRESHOLD SELECTION FOR EXTREME**  
2 **VALUE ANALYSIS VIA ORDERED GOODNESS-OF-FIT**  
3 **TESTS WITH ADJUSTMENT FOR FALSE DISCOVERY**  
4 **RATE**

5 BY BRIAN BADER<sup>\*,†</sup>, JUN YAN<sup>†,§</sup>, AND XUEBIN ZHANG<sup>‡</sup>

6 *KPMG LLP\**, *University of Connecticut†*, and *Environment and Climate*  
7 *Change Canada‡*

8 Threshold selection is a critical issue for extreme value analy-  
9 sis with threshold-based approaches. Under suitable conditions, ex-  
10 ceedances over a high threshold have been shown to follow the gen-  
11 eralized Pareto distribution (GPD) asymptotically. In practice, how-  
12 ever, the threshold must be chosen. If the chosen threshold is too  
13 low, the GPD approximation may not hold and bias can occur. If the  
14 threshold is chosen too high, reduced sample size increases the vari-  
15 ance of parameter estimates. To process batch analyses, commonly  
16 used selection methods such as graphical diagnostics are subjective  
17 and cannot be automated. We develop an efficient technique to evalu-  
18 ate and apply the Anderson–Darling test to the sample of exceedances  
19 above a fixed threshold. In order to automate threshold selection, this  
20 test is used in conjunction with a recently developed stopping rule  
21 that controls the false discovery rate in ordered hypothesis testing.  
22 Previous attempts in this setting do not account for the issue of or-  
23 dered multiple testing. The performance of the method is assessed  
24 in a large scale simulation study that mimics practical return level  
25 estimation. This procedure was repeated at hundreds of sites in the  
26 western US to generate return level maps of extreme precipitation.

27 **1. Introduction.** Extreme value analysis has wide applications in a va-  
28 riety of fields, such as hydrology (e.g., [Katz, Parlange and Naveau, 2002](#))  
29 and climatology (e.g., [Davison and Smith, 1990](#); [Kharin et al., 2013](#)), and  
30 dates back to [Fisher and Tippett \(1928\)](#). A major goal of inferences in these  
31 fields is to estimate the probability of extreme events, often expressed in  
32 terms of return level and return period. A return level with a return period  
33 of  $T = 1/p$  years is a high value that is exceeded with probability  $p$ . In other  
34 words, the average number of events exceeding this level within a  $T$ -year  
35 period is one. Commonly used in extreme value analysis, threshold-based  
36 methods involve modeling data exceeding a suitably chosen high threshold  
37 with the generalized Pareto distribution (GPD) ([Balkema and De Haan,](#)  
38 [1974](#); [Pickands, 1975](#)). Choice of the threshold is critical in obtaining accu-  
39 rate estimates of model parameters and return levels. The threshold should

---

<sup>§</sup>Research partially supported by NSF grant DMS 1521730, University of Connecticut Research Excellence Program, and Environment Canada. The authors thank Prof. Vartan Choulakian for the discussion and insight on approximating the tails of the null distribution of the Anderson–Darling and Cramér–von Mises statistics for generalized Pareto distributed data.

*Keywords and phrases:* batch analysis, exceedance diagnostic, specification test, stopping rule

1 be chosen high enough for the excesses to be well approximated by the GPD  
2 to minimize bias, but not so high to substantially increase the variance of the  
3 estimator due to reduction in the sample size (the number of exceedances).

4 Although it is widely accepted in the statistics community that the peaks-  
5 over-threshold (POT) approach uses data more efficiently than the block  
6 maxima method (e.g., [Caires, 2009](#); [Wang, 1991](#)), it is less utilized in some  
7 fields such as climatology. Even when appropriate data is available for use  
8 with the POT approach, there is a lack of efficient procedures that can  
9 be applied consistently and automatically to select the threshold in each  
10 sample, sometimes numbering in the hundreds or thousands (e.g., [Kharin  
11 et al., 2007, 2013](#)). As a motivating application, consider mapping the an-  
12 nual return levels of daily precipitation for the three west coastal US states  
13 of California, Oregon, and Washington. For the three states, one needs to  
14 repeat the estimation procedure, including threshold selection, at each of  
15 the hundreds of stations. For the whole US, thousands of sites would need  
16 to be processed. A graphical based diagnosis will be difficult to apply consis-  
17 tently across many sites and is clearly impractical. It is desirable to have an  
18 intuitive automated threshold selection procedure to use with POT analysis.

19 Many threshold selection methods are available in the literature; see [Scar-  
20 rott and MacDonald \(2012\)](#), [Caeiro and Gomes \(2015\)](#) and [Langousis et al.  
21 \(2016\)](#) for recent reviews. Graphical diagnosis methods are commonly used  
22 (e.g., [Davison and Smith, 1990](#); [Drees, De Haan and Resnick, 2000](#); [Coles,  
23 2001](#); [Scarrott and MacDonald, 2012](#)), but can be quite subjective and not  
24 appropriate as an automated procedure. Recently, [Northrop, Attalides and  
25 Jonathan \(2017\)](#) take a unique approach by applying Bayesian model av-  
26 eraging to combine inferences from multiple thresholds in order to reduce  
27 the sensitivity in estimates from using a single, fixed threshold. Other selec-  
28 tion methods can be grouped into various categories. One is based on the  
29 asymptotic results about estimators of properties of the tail distribution.  
30 [Langousis et al. \(2016\)](#) detail a few of these procedures for threshold selec-  
31 tion, such as the Jackson ([Jackson, 1967](#)) and Lewis ([Lewis, 1965](#)) kernel  
32 statistics as modified in [Goegebeur, Beirlant and de Wet \(2008\)](#) (based on  
33 the Hill estimator) and an automated version of the mean residual life (MRL)  
34 plot. Computational, resampling-based estimators require sufficient comput-  
35 ing resources, may involve tuning parameters ([Danielsson et al., 2001](#)), and  
36 in some cases is not satisfactory for small samples ([Ferreira, de Haan and  
37 Peng, 2003](#)).

38 A second category of methods are based on goodness-of-fit of the GPD,  
39 where the threshold is selected as the lowest level above which the GPD  
40 provides adequate fit to the exceedances (e.g., [Davison and Smith, 1990](#);

1 Dupuis, 1999; Choulakian and Stephens, 2001; Northrop and Coleman, 2014;  
2 Langousis et al., 2016). Goodness-of-fit tests are simple to understand and  
3 perform, but error control, however, is challenging because of the ordered  
4 nature of the hypotheses, and the usual methods from multiple testing such  
5 as false discovery rate (FDR) (e.g., Benjamini, 2010a,b) cannot be directly  
6 applied. This has not been addressed to the best of our knowledge. Methods  
7 in the third category are based on mixtures of a GPD for the tail and  
8 another distribution for the “bulk” joined at the threshold (e.g., MacDonald  
9 et al., 2011; Wadsworth and Tawn, 2012; Naveau et al., 2016). Treating the  
10 threshold as a parameter to estimate, these methods can account for the  
11 uncertainty from threshold selection in inferences. However, care is needed  
12 to ensure that the bulk and tail models are robust to one another in the case  
13 of misspecification.

14 The simple naive method is *a priori* or fixed threshold selection based on  
15 expertise on the subject matter at hand. Various rules of thumb have been  
16 suggested; for example, selecting the top 10% of the data (e.g., DuMouchel,  
17 1983), or the top square root of the sample size (e.g., Ferreira, de Haan  
18 and Peng, 2003). Such one rule for all is not ideal in climate applications  
19 where high heterogeneity in data properties is the norm. The proportion of  
20 the number of rain days can be very different from wet tropical climates to  
21 dry subtropical climates; therefore the number of exceedances over the same  
22 time period can be very different across different climates. Additionally,  
23 the probability distribution of daily precipitation can also be different in  
24 different climates, affecting the speed of the tail convergence to the GPD  
25 (Raoult and Worms, 2003).

26 We propose an automated threshold selection procedure based on a se-  
27 quence of goodness-of-fit tests with error control for ordered, multiple test-  
28 ing. The recently developed stopping rules in G’Sell et al. (2015), which  
29 control the FDR (Benjamini and Hochberg, 1995; Benjamini and Yekutieli,  
30 2001) for ordered hypotheses, are adapted for use in this setting. They are  
31 applied to the Anderson–Darling (AD) goodness-of-fit test at each candi-  
32 date threshold sequentially from low to high. The application is challeng-  
33 ing in that the asymptotic null distribution of the testing statistic is un-  
34 wieldy (Choulakian and Stephens, 2001), and that parametric bootstrap  
35 puts bounds on the approximate p-values which can be shown to reduce  
36 power of the stopping rules. We propose a fast approximation of the p-value  
37 based on the results of Choulakian and Stephens (2001) to facilitate the  
38 application. The performance of the procedures are investigated in a large  
39 scale simulation study, and recommendations are made. The procedure is  
40 applied to precipitation return level mapping of three west coastal states

1 of the US. Interesting findings are revealed from different stopping rules.  
 2 The automated threshold selection procedure has applications in various  
 3 fields, especially when a consistent application for batch processing of mas-  
 4 sive datasets and of many sites is needed.

5 The outline of the paper is as follows. Section 2 presents the generalized  
 6 Pareto model, its theoretical justification, and how to apply the automated  
 7 sequential threshold testing procedure. Section 3 introduces the tests pro-  
 8 posed to be used in the automated testing procedure. A simulation study  
 9 demonstrates the power of the tests for a fixed threshold under various mis-  
 10 specification settings and it is found that the AD test is most powerful in the  
 11 vast majority of cases. A large scale simulation study in Section 4 demon-  
 12 strates performance of the stopping rules for multiple ordered hypotheses,  
 13 under a plausible misspecified distribution and is compared with competing  
 14 methods. In Section 5, we return to our motivating application and derive re-  
 15 turn levels for extreme precipitation at hundreds of west coastal US stations  
 16 to demonstrate the usage of our method and some practical considerations.  
 17 A final discussion is delivered in Section 6.

18 **2. Automated Sequential Testing Procedure.** Threshold methods  
 19 for extreme value analysis are based on that, under general regularity con-  
 20 ditions, the only possible non-degenerate limiting distribution of properly  
 21 rescaled exceedances of a threshold  $u$  is the GPD as  $u \rightarrow \infty$  (e.g., [Pickands,](#)  
 22 [1975](#)). The  $\text{GPD}(\sigma_u, \xi)$  has cumulative distribution function

$$(1) \quad F(y|\theta) = \begin{cases} 1 - \left(1 + \frac{\xi y}{\sigma_u}\right)^{-1/\xi} & \xi \neq 0, \quad y > 0, \quad 1 + \frac{\xi y}{\sigma_u} > 0, \\ 1 - \exp\left(-\frac{y}{\sigma_u}\right) & \xi = 0, \quad y > 0, \end{cases}$$

23 where  $\theta = (\sigma_u, \xi)$ ,  $\xi$  is a shape parameter, and  $\sigma_u > 0$  is a threshold-  
 24 dependent scale parameter. The GPD also has the property that for some  
 25 threshold  $v > u$ , the excesses follow a GPD with the same shape parameter,  
 26 but a modified scale  $\sigma_v = \sigma_u + \xi(v - u)$ . For the remainder of the text, the  
 27 subscript  $u$  associated with  $\sigma$  is dropped for ease of presentation.

28 Let  $X_1, \dots, X_n$  be a random sample of size  $n$ . If  $u$  is sufficiently high, the  
 29 exceedances  $Y_i = X_i - u$  for all  $i$  such that  $X_i > u$  are approximately a  
 30 random sample from a GPD. The question is to find the lowest threshold  
 31 such that the GPD fits the sample of exceedances over this threshold ade-  
 32 quately. Our solution is to combine sequential goodness-of-fit testing with  
 33 adjustments for multiplicity in the ordered setting. The first component is  
 34 similar to the approach taken in [Choulakian and Stephens \(2001\)](#) or model  
 35 specification testing (e.g., [Northrop and Coleman, 2014](#); [Wadsworth, 2016](#))

1 for the GPD to the exceedances over each candidate threshold in an increas-  
 2 ing order. The second component, multiple testing in the special ordered  
 3 setting, is handled by the stopping rules in G'Sell et al. (2015).

4 Consider a fixed set of candidate thresholds  $u_1 < \dots < u_l$ . For each  
 5 threshold, there will be  $n_i$  excesses,  $i = 1, \dots, l$ . The sequence of null hy-  
 6 potheses can be stated as

7  $H_0^{(i)}$ : The distribution of the  $n_i$  exceedances above  $u_i$  follows the GPD.

8 For a fixed  $u_i$ , many tests are available for this  $H_0^{(i)}$ . An automated pro-  
 9 cedure can begin with  $u_1$  and continue until some threshold  $u_i$  provides an  
 10 acceptance of  $H_0^{(i)}$  (Choulakian and Stephens, 2001; Thompson et al., 2009).  
 11 The problem, however, is that unless the test has high power, an acceptance  
 12 may happen at a low threshold by chance and, thus, the data above the  
 13 chosen threshold is contaminated. One could also begin at the threshold  $u_l$   
 14 and descend until a rejection occurs, but this would result in an increased  
 15 type I error rate. The multiple testing problem obviously needs to be ad-  
 16 dressed, and the issue here is especially challenging because these tests are  
 17 ordered; if  $H_0^{(i)}$  is rejected, then  $H_0^{(k)}$  has been rejected for all  $1 \leq k < i$ .  
 18 Despite the extensive literature on multiple testing and the more recent de-  
 19 velopments on FDR control and its variants (e.g., Benjamini and Hochberg,  
 20 1995; Benjamini and Yekutieli, 2001; Benjamini, 2010a,b), no definitive pro-  
 21 cedure has been available for error control in ordered tests until the recent  
 22 work of G'Sell et al. (2015).

23 We adapt the ForwardStop rule of G'Sell et al. (2015) to the sequential  
 24 testing of (ordered) null hypotheses  $H_1, \dots, H_l$ . Let  $p_1, \dots, p_l \in [0, 1]$  be the  
 25 corresponding p-values of the  $l$  hypotheses. G'Sell et al. (2015) transform  
 26 the sequence of p-values to a monotone sequence and then apply the original  
 27 method of Benjamini and Hochberg (1995) on the monotone sequence. The  
 28 rejection rule is constructed by returning a cutoff  $\hat{k}$  such that  $H_1, \dots, H_{\hat{k}}$  are  
 29 rejected. If no  $\hat{k} \in \{1, \dots, l\}$  exists, then no rejection is made. ForwardStop  
 30 is given by

$$(2) \quad \hat{k}_F = \max \left\{ k \in \{1, \dots, l\} : -\frac{1}{k} \sum_{i=1}^k \log(1 - p_i) \leq \alpha \right\},$$

31 where  $\alpha$  is a pre-specified level.

32 Under the assumption of independence among the tests, ForwardStop was  
 33 shown to control the FDR at level  $\alpha$ . In our setting, stopping at  $k$  implies  
 34 that goodness-of-fit of the GPD to the exceedances at the first  $k$  thresholds  
 35  $\{u_1, \dots, u_k\}$  is rejected. In other words, the first set of  $k$  null hypotheses

1  $\{H_1, \dots, H_k\}$  is rejected. At each  $H_0^{(i)}$ , ForwardStop is a transformed average  
 2 of the previous and current p-values.

3 Another stopping rule developed in G'Sell et al. (2015), StrongStop, could  
 4 be applied in the same manner. However, it provides stronger error control  
 5 than ForwardStop; it controls the familywise error rate instead of FDR. In  
 6 that regard, StrongStop is less desirable for this application as the stricter  
 7 error control leads to decreased power-to-reject and generally results in se-  
 8 lected thresholds that are too low.

9 ForwardStop, combined with the sequential hypothesis testing, provides  
 10 an automated selection procedure — all that is needed is the level of desired  
 11 error control and a set of thresholds. A caveat is that the p-values of the  
 12 sequential tests here are dependent, unlike the setup of G'Sell et al. (2015).  
 13 Nonetheless, ForwardStop may still provide some reasonable error control as  
 14 its counter part in the non-sequential multiple testing scenario (Benjamini  
 15 and Yekutieli, 2001; Blanchard and Roquain, 2009). A simulation study is  
 16 carried out in Section 4 to assess the empirical properties.

17 Additionally, it is of worth to note that there are two potential factors in  
 18 the application to real data that may directly affect threshold selection using  
 19 the ForwardStop procedure. First, by using critical values of goodness-of-  
 20 fit tests intended for continuous data, applied to quantized data, resulting  
 21 p-values have the potential to be underestimated (e.g., Langousis et al.,  
 22 2016). Second, ignoring serial dependence in the excesses can lead to spurious  
 23 precision in inferences about the GP parameters. Although ignoring such  
 24 dependence does not introduce bias, it does underestimate standard errors  
 25 (Fawcett and Walshaw, 2007), which in turn can influence the variability in  
 26 goodness-of-fit test critical values.

27 **3. The Tests.** The automated procedure can be applied with any valid  
 28 test for each hypothesis  $H_0^{(i)}$  corresponding to threshold  $u_i$ . Four existing  
 29 goodness-of-fit tests that can be used are presented. Because the stopping  
 30 rules are based on transformed p-values, it is desirable to have testing statis-  
 31 tics whose p-values can be accurately measured; bootstrap based tests that  
 32 put a lower bound on the p-values (1 divided by the bootstrap sample size)  
 33 may lead to premature stopping. For the remainder of this section, the su-  
 34 perscript  $i$  is dropped. We consider the goodness-of-fit of GPD to a sample  
 35 of size  $n$  of exceedances  $Y = X - u$  above a fixed threshold  $u$ .

36 **3.1. Anderson–Darling and Cramér–von Mises Tests.** The AD and the  
 37 Cramér–von Mises tests for the GPD have been studied in detail (Choulakian  
 38 and Stephens, 2001). Let  $\hat{\theta}_n$  be the maximum likelihood estimator (MLE)  
 39 of  $\theta$  under  $H_0$  from the the observed exceedances. Make the probability

1 integral transformation based on  $\hat{\theta}_n$   $z_{(i)} = F(y_{(i)}|\hat{\theta}_n)$ , as in (1), for the order  
 2 statistics of the exceedances  $y_{(1)} < \dots < y_{(n)}$ . The AD statistic is

$$A_n^2 = -n - \frac{1}{n} \sum_{i=1}^n (2i-1) \left[ \log(z_{(i)}) + \log(1 - z_{(n+1-i)}) \right].$$

3 The Cramér–von Mises statistic is

$$W_n^2 = \sum_{i=1}^n \left[ z_{(i)} - \frac{2i-1}{2n} \right]^2 + \frac{1}{12n}.$$

4 The AD statistic is a slight modification of the Cramér–von Mises statistic,  
 5 with more weight to observations in the tail of the distribution.

6 The asymptotic distributions of  $A_n^2$  and  $W_n^2$  are unwieldy, both being  
 7 sum of weighted chi-squared variables with one degree of freedom with  
 8 weight found from the eigenvalues of an integral equation (Choulakian and  
 9 Stephens, 2001, Section 6). The distributions depend only on the estimate  
 10 of  $\xi$ . The tests are often applied by referring to a table of a few upper tail  
 11 percentiles of the asymptotic distributions (Choulakian and Stephens, 2001,  
 12 Table 2), or through bootstrap. In either case, the p-values are truncated  
 13 by a lower bound. Such truncation of a smaller p-value to a larger one can  
 14 be proven to weaken the ForwardStop rule given in (2). In order to apply  
 15 these tests in the automated sequential setting, more accurate p-values for  
 16 the tests are needed.

17 We provide two remedies to the table in Choulakian and Stephens (2001).  
 18 First, for  $\xi$  values in the range of  $(-0.5, 1)$ , which is applicable for most  
 19 applications, we enlarge the table to a much finer resolution through a pre-  
 20 set Monte Carlo method. For each  $\xi$  value from  $-0.5$  to  $1$  incremented by  
 21  $0.01$ , 2,000,000 replicates of  $A_n^2$  and  $W_n^2$  are generated with sample size  
 22  $n = 1,000$  to approximate their asymptotic distributions. A grid of upper  
 23 percentiles from  $0.999$  to  $0.001$  for each  $\xi$  value is produced and saved in a  
 24 table for fast future reference. Therefore, if estimate  $\hat{\xi}_n$  and the test statistic  
 25 falls in the range of the table, the p-value is computed via interpolation  
 26 on the logarithmic scale. That is, the p-value for an observed test statistic  
 27 is found by taking the distance (on the log-scale) between the two nearest  
 28 critical values found in the table, and scaling the distances to a maximum  
 29 range of  $0.001$ . For example, if the test statistic falls at the one-quarter  
 30 mark between two critical values (on the log-scale), its p-value is equal to  
 31 the smaller critical value's p-value +  $0.00025$ .

32 The second remedy is for observed test statistics that are greater than that  
 33 found in the table (implied p-value less than  $0.001$ ). As Choulakian pointed

1 out (in a personal communication), the tails of the asymptotic distributions  
 2 are exponential, which can be confirmed using the available tail values in the  
 3 table. For a given  $\hat{\xi}_n$ , regressing  $-\log(\text{p-value})$  on the upper tail percentiles  
 4 in the table, for example, from 0.05 to 0.001, gives a linear model that can  
 5 be extrapolated to approximate the p-value of observed statistics outside of  
 6 the range of the table. This approximation of extremely small p-values help  
 7 reduce loss of power in the stopping rules.

8 The two remedies make the two tests very fast and are applicable for most  
 9 applications with  $\xi \in (-0.5, 1)$ . For  $\xi$  values outside of  $(-0.5, 1)$ , although  
 10 slow, one can use parametric bootstrap to obtain the distribution of the test  
 11 statistic, understanding that the p-value has a lower bound. The methods  
 12 are implemented in R package `eva` (Bader and Yan, 2015).

13 *3.2. Moran's Test and Rao's Score Test.* Two other tests are compared  
 14 for reference. Moran's goodness-of-fit test is a byproduct of the maximum  
 15 product spacing (MPS) estimation for estimating the GPD parameters. The  
 16 maximized objective function (i.e., evaluated at the MPS estimators) is  
 17 Moran's statistic (Moran, 1953). Cheng and Stephens (1989) showed that,  
 18 under the GPD null hypothesis, Moran's statistic when properly centered  
 19 and scaled, has an asymptotic chi-square approximation. Wong and Li (2006)  
 20 show empirically the test holds its size for samples as small as ten.

21 Northrop and Coleman (2014) considered a piecewise (varying) represen-  
 22 tation of the shape parameter of the GPD at intervals between a set of  
 23 thresholds. Under  $H_0$  it is assumed that the shape parameter holds a con-  
 24 stant value across all intervals. Tests of this hypothesis are developed based  
 25 on a multiple-threshold penultimate model, using a construction of Rao's  
 26 score test.

27 *3.3. A Power Study.* The power of the four goodness-of-fit tests were  
 28 examined in an individual, non-sequential testing framework. The data gen-  
 29 erating schemes in Choulakian and Stephens (2001) were used, some of which  
 30 are very difficult to distinguish from the GPD:

- 31 • Gamma with shape 2 and scale 1.
- 32 • Standard lognormal distribution (mean 0 and scale 1 on log scale).
- 33 • Weibull with scale 1 and shape 0.75.
- 34 • Weibull with scale 1 and shape 1.25.
- 35 • 50/50 mixture of GPD(1, -0.4) and GPD(1, 0.4).
- 36 • 50/50 mixture of GPD(1, 0) and GPD(1, 0.4).
- 37 • 50/50 mixture of GPD(1, -0.25) and GPD(1, 0.25).

38 Finally, the GPD(1, 0.25) was also used to check the type I error rate. Four



1 sample sizes were considered: 50, 100, 200, 400. For each scenario, 10,000  
 2 samples are generated. The four tests were applied to each sample, with a  
 3 rejection recorded if the p-value is below 0.05. An additional requirement  
 4 of the score test is to select a set of intervals to constitute the piecewise  
 5 representation of the shape parameter (Section 1.2, [Northrop and Coleman, 2014](#));  
 6 we set these according to the deciles of the generated data. The  
 7 likelihood under each specification was maximized with the default setting  
 8 of R function `optim`, the Nelder–Mead method.

9 [Table 1 about here.]

10 The rejection rates are summarized in Table 1. Samples in which the MLE  
 11 failed were removed, which accounts for roughly 10.8% of the Weibull sam-  
 12 ples with shape 1.25 and sample size 400, and around 10.7% for the Gamma  
 13 distribution with sample size 400. Decreasing the sample size in these cases  
 14 actually decreases the percentage of failed MLE samples. This may be due  
 15 to the shape of these two distributions, which progressively become more  
 16 distinct from the GPD as their shape parameters increase. In particular,  
 17 the Weibull changes shape dramatically as it crosses 1, which is the case  
 18 here (0.75 to 1.25). In the other distribution cases, no setting resulted in  
 19 more than a 0.3% failure rate. As expected, all tests appear to hold their  
 20 sizes, and their powers all increase with sample size. The mixture of two  
 21 GPDs is the hardest to detect. For the GPD mixture of shape parameters 0  
 22 and 0.4, quantile matching between a single large sample of generated data  
 23 and the fitted GP distribution shows a high degree of similarity. In the vast  
 24 majority of cases, the AD test appears to have the highest power, followed  
 25 by the Cramér–von Mises test. The R-language code used to perform this  
 26 power study can be found in the Supplementary Materials ([Bader, Yan and Zhang, 2017](#)).  
 27

28 **4. Simulation Study of the Automated Procedure.** It is of inter-  
 29 est to investigate the performance of ForwardStop versus several competing  
 30 methods under misspecification. Empirical work on rainfall extremes (e.g.,  
 31 [Papalexiou and Koutsoyiannis, 2013](#); [Serinaldi and Kilsby, 2014](#)) finds that  
 32 the right tail for daily rainfall is better described by heavy tail distribu-  
 33 tions than exponential types. Following the study of [Roth, Jongbloed and Buishand \(2016\)](#),  
 34 data is generated from a distribution characterized by  
 35 a hazard function which transits smoothly from the hazard function of a  
 36 Weibull distribution to that of a GPD, with the transition occurring near  
 37 some threshold  $u$ . Let  $h_1(x) = \kappa\beta^{-\kappa}x^{\kappa-1}$ ,  $x > 0$ , be the hazard function of  
 38 a Weibull distribution with parameters  $(\kappa, \beta)$ . Let  $h_2(x) = (\sigma + \xi(x - u))^{-1}$ ,  
 39  $x \geq u$  be the hazard function of a GPD with parameters  $(\xi, \sigma)$ . A smooth

1 transition from  $h_1(\cdot)$  to  $h_2(\cdot)$ , with a transition period of length between  
 2  $(u, u + \epsilon)$ , defines a new hazard function

$$h(x) := h_1(x)\eta\left(\frac{x-u}{\epsilon}\right) + h_2(x)\left(1 - \eta\left(\frac{x-u}{\epsilon}\right)\right),$$

3 where

$$\eta(x) = \begin{cases} 1, & x \leq 0, \\ 2x^3 - 3x^2 + 1, & 0 < x < 1, \\ 0, & x \geq 1, \end{cases}$$

4 to ensure that  $h$  has a continuous derivative. The distribution function of  
 5 the data generating mechanism is then

$$F(x) = 1 - \exp(-H(x)) = 1 - \exp\left(-\int_0^x h(z)dz\right).$$

6 More detailed properties of this distribution are discussed in [Roth, Jong-](#)  
 7 [bloed and Buishand \(2016, Section 5.1\)](#); a similar construction consisting of  
 8 a GPD mixture distribution was proposed by [Holden and Haug \(2009\)](#). Four  
 9 different thresholds are considered —  $(u_1, u_2, u_3, u_4) = (17.70, 11.00, 5.91, 3.64)$   
 10 with  $(\sigma_{u_1}, \sigma_{u_2}, \sigma_{u_3}, \sigma_{u_4}) = (11.51, 10.45, 6.80, 4.55)$ . All four distributions  
 11 share the same transition period length  $\epsilon = 0.5$ , GPD parameter  $\xi = 0.15$   
 12 and Weibull parameters  $(\kappa, \beta) = (0.45, 1)$ . Starting from  $u + \epsilon$ , the tail of the  
 13 distribution is completely the tail of the GPD. Our selections correspond to  
 14 the 97.5th, 95th, 90th, and 85th percentiles, respectively. The distribution  
 15 and density functions can be seen in [Figure 1](#). For each setting,  $B = 1, 500$   
 16 samples were generated from this distribution with  $n = 4, 600$ , which can be  
 17 thought of as 50 years of seasonal daily observations (92 days each season).

18 [Fig 1 about here.]

19 The main quantity of interest is the  $N$ -year return level, defined for the  
 20 GPD (e.g., [Coles, 2001, Section 4.3.3](#)) as

$$(3) \quad z_N = \begin{cases} u + \frac{\sigma}{\xi}[(Nn_y\zeta_u)^\xi - 1], & \xi \neq 0, \\ u + \sigma \log(Nn_y\zeta_u), & \xi = 0, \end{cases}$$

21 for a given threshold  $u$ , where  $n_y$  is the number of observations per year,  
 22 and  $\zeta_u$  is the rate, or proportion of the data exceeding  $u$ . In precise terms,  
 23 this is a high quantile of the GPD, which has the interpretation as the level  
 24 expected to be exceeded in a single year with probability  $1/N$ . In estimation,  
 25  $\zeta_u$  is set to the proportion of threshold exceedances.

1 In addition to ForwardStop, three competing methods in return level esti-  
 2 mation were used for comparison. The first was the so-called ‘rules-of-thumb’  
 3 (e.g., [Ferreira, de Haan and Peng, 2003](#); [DuMouchel, 1983](#)), which simply  
 4 uses a fixed (upper) fraction of the data to fit the GPD. In this setting, the  
 5 top 15%, 10%, 5%, and square root (1.5%) of the sample size were used in  
 6 each simulation to fit the GPD. The second, an automated procedure based  
 7 on the MRL plot using weighted least squares (WLS) and detailed in [Lan-](#)  
 8 [gousis et al. \(2016, Section 2.2\)](#), chooses the threshold that minimizes the  
 9 MSE of the WLS fit. The third was the unadjusted alternative to Forward-  
 10 Stop. It can be implemented two ways — both relying on the raw p-values  
 11 from each threshold tested and proceeding sequentially. Specifically, the first  
 12 version (denoted as Raw Up) begins at the lowest threshold and selects the  
 13 first (lowest) threshold which is accepted. If all are rejected, the maximum  
 14 threshold is selected. The second (denoted as Raw Down) starts from the  
 15 largest threshold and proceeds until a rejection of the test occurs; then the  
 16 threshold before the rejection is selected. If a rejection occurs on the first  
 17 (highest) threshold, that threshold is used. ForwardStop has the same di-  
 18 rection of operation as Raw Up, so, it also chooses the largest threshold if  
 19 all are rejected.

20 The estimated 50-year return levels were compared using the chosen  
 21 threshold. For each sample, two sets of percentiles were used to generate  
 22 thresholds to test. The first takes 10 percentiles, in increments of 5 be-  
 23 ginning at 50 and ending at 95, the second takes 20 percentiles, in incre-  
 24 ments of 2.5 beginning at 50 and ending at 97.5. Three significance levels  
 25  $\alpha \in \{0.05, 0.10, 0.20\}$  were used for testing with ForwardStop and the two  
 26 unadjusted procedures. For the automated MRL plot method, a threshold  
 27 was selected for each set of thresholds. For a given return period  $N$ , the root  
 28 mean squared error (RMSE) of an estimator is calculated as

$$\sqrt{\frac{1}{B} \sum_{i=1}^B (\hat{z}_N^i - z_N)^2}.$$

29 Figure 2 shows the RMSE comparison of all the methods considered; the  
 30 corresponding plot for bias and variance (standard deviation) are available  
 31 in the Supplementary Materials ([Bader, Yan and Zhang, 2017](#)).

32 [Fig 2 about here.]

33 As expected, for the ‘top’ methods with a fixed threshold, the resulting  
 34 error in estimation is highly dependent on the chosen fixed threshold. The  
 35 correct fixed threshold fixed method always has smallest RMSE. However,

1 when the threshold is fixed at an incorrect percentile, it is outperformed  
2 by ForwardStop and the unadjusted Raw Up/Down procedures. The MRL  
3 method has varying performance, dependent on the exceedance probability.  
4 Another property of the MRL is that a threshold is always selected and  
5 thus the steps for further diagnostic testing is unclear. The absolute bias  
6 for ForwardStop is generally lower than Raw Up/Down. Compared with  
7 ForwardStop, Raw Down performs better for small exceedance probabilities  
8 and worse for larger values (variance increases); Raw Up performs better for  
9 larger exceedance probabilities and worse for smaller values (bias increases).  
10 ForwardStop hedges against both cases and provides similar or smaller error  
11 than the better of Raw Up and Raw Down. This is an intuitive result — as  
12 the exceedance probability decreases, it is more likely for Raw Up to stop too  
13 early (increased bias). Similarly, as the exceedance probability increases, it is  
14 more likely for Raw Down to stop too high (increased variance). ForwardStop  
15 guards against either of these cases.

16 An alternative simulation study is reported in the Supplementary Mate-  
17 rials (Bader, Yan and Zhang, 2017).

18 **5. Return Level Mapping of Extreme Precipitation.** Hydrolo-  
19 gists are often interested in estimating return levels of extreme precipitation  
20 (Katz, Parlange and Naveau, 2002). These return levels are provided as maps  
21 for infrastructure design for risk management and land planning (Blanchet  
22 and Lehning, 2010; Lateltin and Bonnard, 1999). Estimating return levels  
23 across many sites consistently is often required to produce such maps. Our  
24 methodology uses automated threshold selection to generate return level  
25 maps through batch processing of data from a large number of sites. Here  
26 we consider a return level map of extreme precipitation in the three west-  
27 ern US coastal states of California, Oregon, and Washington with diverse  
28 climates to demonstrate its application and practical usage. The automated  
29 procedure in Section 2 provides a quick and consistent way to obtain an  
30 accurate map without the need to inspect every site, and without taking a  
31 homogeneous approach.

32 Daily precipitation data is available for tens of thousands of surface sites  
33 around the world via the Global Historical Climatology Network (GHCN).  
34 A description of the data network can be found in Menne et al. (2012). After  
35 an initial screening to remove sites with less than 50 years of available data,  
36 there were 720 remaining sites across the three chosen coastal states. As the  
37 annual maximum daily amount of precipitation mainly occurs in winter, only  
38 the winter season (November to March) observations were used in modeling.

39 A major issue in the analysis of US precipitation data is the sensitivity of

1 statistical procedures in response to quantization. The United States precip-  
 2 itation data has a unit in millimeters (mm) and the data are recorded to the  
 3 nearest hundredth of an inch. As a result, a rough rounding occurs with 0.2  
 4 and 0.3 mm intervals. The inflation effect of quantization on goodness-of-  
 5 fit tests for the GPD (in particular, AD) has been well-known (e.g., [Deidda](#)  
 6 [and Puliga, 2006, 2009](#); [Langousis et al., 2016](#)). Quantization pushes the null  
 7 distribution of the AD statistic to the right; the p-value obtained by posi-  
 8 tioning the observed statistic with quantized data to the null distribution  
 9 from continuous data is smaller than it should be. Rougher rounding means  
 10 bigger change in magnitude of the null distribution. Given a rounding level,  
 11 the smaller the scale parameter of the null GPD, the bigger the increase in  
 12 the magnitude of the AD test statistic. To accommodate the quantization  
 13 issue, we seek to adjust the test statistic instead of the sampling distribu-  
 14 tion, because the sampling distribution depends on the scale parameters and  
 15 needs to be approximated by Monte Carlo.

16 We approached the quantization issue by treating the observed data as  
 17 uniformly interval censored. That is, for an observed value of  $x_i > 0$ , the ac-  
 18 tual observation is treated as interval censored by the interval  $[x_i - \delta/2, x_i +$   
 19  $\delta/2)$ , where  $\delta = 0.254$ . To calculate the AD statistic, we perturbed each  
 20 observed value  $x_i$  by a random draw from the uniform distribution over  
 21  $[-\delta/2, \delta/2)$ . We replicated the jittering process  $K$  times, and took the me-  
 22 dian of the resulting  $K$  AD statistics as our testing statistic. Its p-value was  
 23 then found from the null distribution obtained for continuous data via the  
 24 procedure in Section 3.1. As observed in the simulation study reported in  
 25 the Supplementary Materials ([Bader, Yan and Zhang, 2017](#)), this jittering  
 26 brings the AD statistic from quantized data much closer to the continuous  
 27 version than the non-jittered version. The distribution of the jittered version  
 28 of the AD statistic is approximated reasonably well by the distribution of  
 29 the AD statistic for continuous data under varying degrees of quantization  
 30 when  $\sigma_u$  is ‘large’. Based on the experiments, sites are screened by taking  
 31  $\sigma_{u=F^{-1}(0.7)} > 5.3$ , which is the 25th percentile of estimated scale parameters  
 32 estimated at all sites in an exploratory data analysis using the top 30% of  
 33 the data at each site. This results in 501 sites to be used. A fully satisfactory  
 34 solution to handle the quantization is noted as a topic of future work for its  
 35 practical importance.

36 Another potential source of uncertainty in the goodness-of-fit testing pro-  
 37 cedure is the presence of serial dependence in the excesses. The goodness-  
 38 of-fit tests assume independence in the excesses under the null hypothesis  
 39 and thus significant departures may affect the testing conclusion. One way  
 40 to check for this is the extremal index ([Leadbetter et al., 1989](#)), which is a

1 measure of the clustering of the underlying process at extreme levels and  
2 quantifies such dependence. It can take values from 0 to 1, with independent  
3 series exhibiting a value of exactly 1. To get a sense for the properties of  
4 series in this dataset, the extremal index was estimated using the so-called  
5 intervals estimator of [Ferro and Segers \(2003\)](#) for each of the 501 selected  
6 sites at the 70th percentile threshold via the R package `texmex` ([Southworth  
7 and Heffernan, 2012](#)). In summary, 199 sites (40%) have an estimated ex-  
8 tremal index of 1, with a minimum value of 0.819. So the serial dependence  
9 was not considered largely influential.

10 Candidate thresholds were chosen based on the data percentiles. For each  
11 site, the set of thresholds was formed by taking the 70th to 98th percentiles  
12 in increments of 2, resulting in 15 thresholds to test. The 98th percentile  
13 is chosen as the upper limit to ensure a sufficient amount of data is avail-  
14 able for parameter estimation. The FS rule is applied using significance level  
15  $\alpha = 0.05$ . If all thresholds were rejected at a site, the 98th percentile was  
16 used to estimate the GPD parameters and corresponding return levels for  
17 that site. Out of the 501 sites, 33 (7%) had all thresholds rejected. This is  
18 an indication that these cases need to be looked at more thoroughly. One  
19 possible explanation is that the quantized data led to over rejection in the  
20 AD tests, and, hence, the threshold selected tends to be over high. If an in-  
21 direct solution is desired, one may relax the distributional assumptions (e.g.,  
22 [Papastathopoulos and Tawn, 2013](#); [Nadarajah and Eljabri, 2013](#)) and follow  
23 the same procedure. It is of worth to compare the selected thresholds at  
24 each site found by various methods. Pairwise plots of selected thresholds for  
25 the FS, unadjusted, and MRL methods are presented in the Supplementary  
26 Materials ([Bader, Yan and Zhang, 2017](#)).

27 [Fig 3 about here.]

28 With the automatically selected thresholds for the 501 sites, the return  
29 levels were estimated. Figure 3 shows map of the 50-, 100-, and 250-year es-  
30 timated return levels of the three states. Note that only a very small number  
31 of sites in eastern Washington/Oregon and the Great Valley of California  
32 show on the map, with low return levels. The screening process removed  
33 many sites from these specific geographic areas, which had a smaller per-  
34 centage of precipitation days than the remaining 501 analyzed on average  
35 (27% versus 38%) in the winter season.. These areas are known dry areas  
36 because the moist air from the Pacific has already lost much of its moisture  
37 on the windward west side of the mountains, and as it descends on the lee-  
38 ward east side of the mountains, it warms adiabatically, making it less likely  
39 for the air to saturate and form precipitation. The selective feature of the

1 FS procedure even after the screening process is a desirable as it suggests  
2 not to fit GPD at even the highest threshold at these sites, a guard that is  
3 not available from those unconditionally applied, one-for-all rules.

4 **6. Discussion.** We propose an intuitive and comprehensive methodol-  
5 ogy for automated threshold selection in the peaks over threshold approach.  
6 In addition, it is less computationally intensive than some competing resam-  
7 pling or bootstrap procedures (Danielsson et al., 2001; Ferreira, de Haan and  
8 Peng, 2003). Automation and efficiency is required when prediction using  
9 the peaks-over-threshold approach is desired at a large number of sites. This  
10 is achieved through sequentially testing a set of thresholds for goodness-of-  
11 fit to the GPD. This general methodology has been applied previously (e.g.,  
12 Choulakian and Stephens, 2001; Thompson et al., 2009); however these did  
13 not account for the multiplicity issue. That is, they selected a threshold  
14 by checking raw p-values against the desired significance level until a rejec-  
15 tion occurred. We apply the recently developed stopping rule ForwardStop  
16 (G'Sell et al., 2015), which transforms the results of ordered, sequentially  
17 tested hypotheses to control the false discovery rate. There is a slight caveat  
18 in our setting, that the tests are not independent, but it can be demon-  
19 strated via simulation that the stopping rules in G'Sell et al. (2015) still  
20 provide reasonable error control here.

21 Four tests are compared in terms of power to detect departures from the  
22 GPD at a single, fixed threshold and it is found that the AD test has the  
23 most power in various non-null settings. Choulakian and Stephens (2001)  
24 derived the asymptotic null distribution of the AD test statistic. However  
25 this requires solving an integral equation. Our contribution, with some ad-  
26 vice from Professor Choulakian, provides an approximate, but accurate and  
27 computationally efficient version of this test. To investigate the performance  
28 of ForwardStop in conjunction with the AD test, a large scale simulation  
29 study was conducted. Data is generated from a plausible distribution – mis-  
30 specified below a certain threshold and generated from the null GPD above,  
31 smoothed via a transition function. In each replicate, the squared error is  
32 recorded for various parameters using ForwardStop and compared with com-  
33 peting alternative procedures. The results of this simulation are mixed – no  
34 procedure outperformed in all settings. ForwardStop is, however, less sen-  
35 sitive to the starting threshold than the unadjusted versions of sequential  
36 hypothesis testing (Raw Up and Down) and it is not forced to always select  
37 a threshold, unlike the MRL approach.

38 The entire methodology using ForwardStop is applied to daily precipi-  
39 tation data at hundreds of sites in three U.S. west coast states, with the

1 goal of creating a return level map. However, quantization present in the  
2 U.S. GHCN dataset introduces additional uncertainty and in this instance,  
3 increased conservatism when performing goodness-of-fit testing on the quan-  
4 tized data. To handle this, the data are treated as (uniformly) interval cen-  
5 sored and an adjusted testing statistic is found via realized samples. As is  
6 discussed in the Supplementary Materials (Bader, Yan and Zhang, 2017),  
7 this solution is still conservative (i.e., smaller p-values, higher thresholds  
8 selected). The most precise way is to obtain the null distribution of the  
9 Anderson–Darling statistic under quantization, but this is cumbersome as  
10 the distribution depends on both the relative size of  $\sigma_u$  to the quantization  
11 level and  $\xi$ . Sites across eastern Washington/Oregon were screened out or  
12 estimated to have smaller return levels, which is consistent with the known  
13 climate of that region.

14 Temporal or covariate varying thresholds, as discussed in Roth et al.  
15 (2012) and Northrop and Jonathan (2011), is an obvious extension to this  
16 work. However, one particular complication that arises is performing model  
17 selection (i.e., what covariates to include), while concurrently testing for  
18 goodness-of-fit to various thresholds. It is clear that threshold selection will  
19 be dependent on the choice of model. Another possible extension involves  
20 testing for overall goodness-of-fit across sites (e.g., Roth, Jongbloed and  
21 Buishand, 2016). In this way, a fixed or quantile regression based threshold  
22 may be predetermined and then tested simultaneously across sites. In this  
23 setup both spatial and temporal dependence need to be taken into account.  
24 Handling this requires some care due to censoring (e.g., Dey and Yan, 2015,  
25 Section 2.5.2). In other words, it is not straightforward to capture the tem-  
26 poral dependence as exceedances across sites are not guaranteed to occur at  
27 the same points in time.

28 An important issue that needs to be addressed, both here and the field  
29 of extremes in general is quantization. It has been studied in detail (Dei-  
30 dda and Puliga, 2006, 2009; Langousis et al., 2016). The most obvious way  
31 to handle quantization is via simulation and/or bootstrap approaches. Di-  
32 rect application of this becomes cumbersome when testing the GPD for  
33 goodness-of-fit at a large number of thresholds since the testing statistic un-  
34 der quantization is dependent on both shape and scale parameters. Its effect  
35 has been reduced in this data analysis but future work is needed to directly  
36 accommodate quantization, particularly in the sequential testing setting.

## SUPPLEMENTARY MATERIAL

37 **Additional Simulation Results and Data Analysis**  
38 (doi: COMPLETED BY THE TYPESETTER; .pdf). Material consisting



1 of R code for the power study, additional simulation results, and analysis  
2 related to the application.

### 3 **References.**

- 4 BADER, B. and YAN, J. (2015). *eva: Extreme Value Analysis with Goodness-of-Fit Testing*  
5 R package version 0.1.2.
- 6 BADER, B., YAN, J. and ZHANG, X. (2017). Supplementary Materials to “Automated  
7 Threshold Selection for Extreme Value Analysis via Ordered Goodness-of-Fit Tests  
8 with Adjustment for False Discovery Rate”. *Annals of Applied Statistics*.
- 9 BALKEMA, A. A. and DE HAAN, L. (1974). Residual Life Time at Great Age. *The Annals*  
10 *of Probability* **2** 792–804.
- 11 BENJAMINI, Y. (2010a). Discovering the False Discovery Rate. *Journal of the Royal Sta-*  
12 *tistical Society: Series B (Statistical Methodology)* **72** 405–416.
- 13 BENJAMINI, Y. (2010b). Simultaneous and Selective Inference: Current Successes and Fu-  
14 ture Challenges. *Biometrical Journal* **52** 708–721.
- 15 BENJAMINI, Y. and HOCHBERG, Y. (1995). Controlling the False Discovery Rate: A Practi-  
16 cal and Powerful Approach to Multiple Testing. *Journal of the Royal Statistical Society.*  
17 *Series B* **57** 289–300.
- 18 BENJAMINI, Y. and YEKUTIELI, D. (2001). The Control of the False Discovery Rate in  
19 Multiple Testing under Dependency. *The Annals of Statistics* **29** 1165–1188.
- 20 BLANCHARD, G. and ROQUAIN, É. (2009). Adaptive False Discovery Rate Control under  
21 Independence and Dependence. *The Journal of Machine Learning Research* **10** 2837–  
22 2871.
- 23 BLANCHET, J. and LEHNING, M. (2010). Mapping Snow Depth Return Levels: Smooth  
24 Spatial Modeling versus Station Interpolation. *Hydrology and Earth System Sciences*  
25 **14** 2527–2544.
- 26 CAEIRO, F. and GOMES, M. I. (2015). Threshold Selection in Extreme Value Analysis. In  
27 *Extreme Value Modeling and Risk Analysis: Methods and Applications* (D. K. Dey and  
28 J. Yan, eds.) 69–82. CRC Press.
- 29 CAIRES, S. (2009). A Comparative Simulation Study of the Annual Maxima and the Peaks-  
30 over-Threshold Methods Technical Report, SBW-Belastingen: subproject ‘Statistics’.  
31 Deltares Report 1200264-002.
- 32 CHENG, R. C. H. and STEPHENS, M. A. (1989). A Goodness-of-Fit Test Using Moran’s  
33 Statistic with Estimated Parameters. *Biometrika* **76** 385–392.
- 34 CHOULAKIAN, V. and STEPHENS, M. A. (2001). Goodness-of-Fit Tests for the Generalized  
35 Pareto Distribution. *Technometrics* **43** 478–484.
- 36 COLES, S. (2001). *An Introduction to Statistical Modeling of Extreme Values*, 1 ed.  
37 Springer.
- 38 DANIELSSON, J., DE HAAN, L., PENG, L. and DE VRIES, C. G. (2001). Using a Bootstrap  
39 Method to Choose the Sample Fraction in Tail Index Estimation. *Journal of Multivari-*  
40 *ate Analysis* **76** 226–248.
- 41 DAVISON, A. C. and SMITH, R. L. (1990). Models for Exceedances Over High Thresholds.  
42 *Journal of the Royal Statistical Society. Series B (Methodological)* **52** 393–442.
- 43 DEIDDA, R. and PULIGA, M. (2006). Sensitivity of Goodness-of-Fit Statistics to Rainfall  
44 Data Rounding Off. *Physics and Chemistry of the Earth, Parts A/B/C* **31** 1240–1251.
- 45 DEIDDA, R. and PULIGA, M. (2009). Performances of Some Parameter Estimators of the  
46 Generalized Pareto Distribution over Rounded-Off Samples. *Physics and Chemistry of*  
47 *the Earth, Parts A/B/C* **34** 626–634.

- 1 DEY, D. K. and YAN, J., eds. (2015). *Extreme Value Modeling and Risk Analysis: Methods*  
2 *and Applications*. CRC Press.
- 3 DREES, H., DE HAAN, L. and RESNICK, S. (2000). How to Make a Hill Plot. *The Annals*  
4 *of Statistics* **28** 254–274.
- 5 DUMOUCHEL, W. H. (1983). Estimating the Stable Index  $\alpha$  in Order to Measure Tail  
6 Thickness: A Critique. *The Annals of Statistics* **11** 1019–1031.
- 7 DUPUIS, D. J. (1999). Exceedances over High Thresholds: A Guide to Threshold Selection.  
8 *Extremes* **1** 251–261.
- 9 FAWCETT, L. and WALSHAW, D. (2007). Improved Estimation for Temporally Clustered  
10 Extremes. *Environmetrics* **18** 173–188.
- 11 FERREIRA, A., DE HAAN, L. and PENG, L. (2003). On Optimising the Estimation of High  
12 Quantiles of a Probability Distribution. *Statistics* **37** 401–434.
- 13 FERRO, C. A. and SEGERS, J. (2003). Inference for Clusters of Extreme Values. *Journal*  
14 *of the Royal Statistical Society: Series B (Statistical Methodology)* **65** 545–556.
- 15 FISHER, R. A. and TIPPETT, L. H. C. (1928). Limiting forms of the frequency distribution  
16 of the largest or smallest member of a sample. In *Mathematical Proceedings of the*  
17 *Cambridge Philosophical Society* **24** 180–190. Cambridge Univ Press.
- 18 GOEGEBEUR, Y., BEIRLANT, J. and DE WET, T. (2008). Linking Pareto-tail Kernel  
19 Goodness-of-Fit Statistics with Tail Index at Optimal Threshold and Second Order  
20 Estimation. *Revstat* **6** 51–69.
- 21 G’SSELL, M. G., WAGER, S., CHOULDECHOVA, A. and TIBSHIRANI, R. (2015). Sequential  
22 Selection Procedures and False Discovery Rate Control. *Journal of the Royal Statistical*  
23 *Society: Series B (Statistical Methodology)*. Forthcoming.
- 24 HOLDEN, L. and HAUG, O. (2009). A Multidimensional Mixture Model for Unsupervised  
25 Tail Estimation NR-notat SAMBA/09/09. pp 29.
- 26 JACKSON, O. (1967). An Analysis of Departures from the Exponential Distribution. *Journal*  
27 *of the Royal Statistical Society. Series B (Methodological)* 540–549.
- 28 KATZ, R. W., PARLANGE, M. B. and NAVEAU, P. (2002). Statistics of Extremes in Hy-  
29 drology. *Advances in Water Resources* **25** 1287–1304.
- 30 KHARIN, V. V., ZWIERS, F. W., ZHANG, X. and HEGERL, G. C. (2007). Changes in  
31 Temperature and Precipitation Extremes in the IPCC Ensemble of Global Coupled  
32 Model Simulations. *Journal of Climate* **20** 1419–1444.
- 33 KHARIN, V. V., ZWIERS, F., ZHANG, X. and WEHNER, M. (2013). Changes in Tem-  
34 perature and Precipitation Extremes in the CMIP5 Ensemble. *Climatic Change* **119**  
35 345–357.
- 36 LANGOUSIS, A., MAMALAKIS, A., PULIGA, M. and DEIDDA, R. (2016). Threshold Detec-  
37 tion for the Generalized Pareto Distribution: Review of Representative Methods and  
38 Application to the NOAA NCDC Daily Rainfall Database. *Water Resources Research*  
39 **52** 2659–2681.
- 40 LATELTIN, O. and BONNARD, C. (1999). Hazard Assessment and Land-Use Planning  
41 in Switzerland for Snow Avalanches, Floods and Landslides Technical Report, World  
42 Meteorological Organization.
- 43 LEADBETTER, M. R., WEISSMAN, I., DE HAAN, L. and ROOTZÉN, H. (1989). On Clus-  
44 tering of High Values in Statistically Stationary Series. *Proc. 4th Int. Meet. Statistical*  
45 *Climatology* **16** 217–222.
- 46 LEWIS, P. A. (1965). Some results on tests for Poisson processes. *Biometrika* **52** 67–77.
- 47 MACDONALD, A., SCARROTT, C. J., LEE, D., DARLOW, B., REALE, M. and RUSSELL, G.  
48 (2011). A Flexible Extreme Value Mixture Model. *Computational Statistics & Data*  
49 *Analysis* **55** 2137–2157.
- 50 MENNE, M. J., DURRE, I., VOSE, R. S., GLEASON, B. E. and HOUSTON, T. G. (2012).

- 1 An overview of the global historical climatology network-daily database. *Journal of*  
2 *Atmospheric and Oceanic Technology* **29** 897–910.
- 3 MORAN, P. A. P. (1953). The Random Division of an Interval-Part II. *Journal of the*  
4 *Royal Statistical Society. Series B (Methodological)* **15** 77–80.
- 5 NADARAJAH, S. and ELJABRI, S. (2013). The Kumaraswamy GP distribution. *Journal of*  
6 *Data Science* **11** 739–766.
- 7 NAVEAU, P., HUSER, R., RIBEREAU, P. and HANNART, A. (2016). Modeling Jointly Low,  
8 Moderate, and Heavy Rainfall Intensities Without a Threshold Selection. *Water Re-*  
9 *sources Research* **52** 2753–2769.
- 10 NORTHROP, P. J., ATTALIDES, N. and JONATHAN, P. (2017). Cross-validators extreme  
11 value threshold selection and uncertainty with application to ocean storm severity.  
12 *Journal of the Royal Statistical Society: Series C (Applied Statistics)* **66** 93–120.
- 13 NORTHROP, P. J. and COLEMAN, C. L. (2014). Improved Threshold Diagnostic Plots for  
14 Extreme Value Analyses. *Extremes* **17** 289–303.
- 15 NORTHROP, P. J. and JONATHAN, P. (2011). Threshold modelling of spatially dependent  
16 non-stationary extremes with application to hurricane-induced wave heights. *Environ-*  
17 *metrics* **22** 799–809.
- 18 PAPALEXIOU, S. M. and KOUTSOYIANNIS, D. (2013). Battle of extreme value distributions:  
19 A global survey on extreme daily rainfall. *Water Resources Research* **49** 187–201.
- 20 PAPASTATHOPOULOS, I. and TAWN, J. A. (2013). Extended generalised Pareto models for  
21 tail estimation. *Journal of Statistical Planning and Inference* **143** 131–143.
- 22 PICKANDS, J. III (1975). Statistical Inference Using Extreme Order Statistics. *The Annals*  
23 *of Statistics* **3** 119–131.
- 24 RAOULT, J.-P. and WORMS, R. (2003). Rate of Convergence for the Generalized Pareto  
25 Approximation of the Excesses. *Advances in Applied Probability* **35** 1007–1027.
- 26 ROTH, M., JONGBLOED, G. and BUIHAND, T. A. (2016). Threshold Selection for Regional  
27 Peaks-over-Threshold Data. *Journal of Applied Statistics* **43** 1291–1309.
- 28 ROTH, M., BUIHAND, T. A., JONGBLOED, G., KLEIN TANK, A. M. G. and VAN ZAN-  
29 TEN, J. H. (2012). A Regional Peaks-over-Threshold Model in a Nonstationary Climate.  
30 *Water Resources Research* **48**.
- 31 SCARROTT, C. and MACDONALD, A. (2012). A Review of Extreme Value Threshold Es-  
32 timation and Uncertainty Quantification. *REVSTAT-Statistical Journal* **10** 33–60.
- 33 SERINALDI, F. and KILSBY, C. G. (2014). Rainfall extremes: Toward reconciliation after  
34 the battle of distributions. *Water resources research* **50** 336–352.
- 35 SOUTHWORTH, H. and HEFFERNAN, J. E. (2012). texmex: Threshold Exceedences and  
36 Multivariate Extremes R package version 1.3.
- 37 THOMPSON, P., CAI, Y., REEVE, D. and STANDER, J. (2009). Automated Threshold  
38 Selection Methods for Extreme Wave Analysis. *Coastal Engineering* **56** 1013–1021.
- 39 WADSWORTH, J. L. (2016). Exploiting Structure of Maximum Likelihood Estimators for  
40 Extreme Value Threshold Selection. *Technometrics* **58** 116–126.
- 41 WADSWORTH, J. L. and TAWN, J. A. (2012). Likelihood-Based Procedures for Thresh-  
42 old Diagnostics and Uncertainty in Extreme Value Modelling. *Journal of the Royal*  
43 *Statistical Society: Series B (Statistical Methodology)* **74** 543–567.
- 44 WANG, Q. J. (1991). The POT Model Described by the Generalized Pareto Distribution  
45 with Poisson Arrival Rate. *Journal of Hydrology* **129** 263–280.
- 46 WONG, T. S. T. and LI, W. K. (2006). A Note on the Estimation of Extreme Value  
47 Distributions Using Maximum Product of Spacings. In *Time Series and Related Topics*  
48 272–283. Institute of Mathematical Statistics.

1 BRIAN BADER  
KPMG LLP  
560 LEXINGTON AVENUE  
NEW YORK, NY 10022  
E-MAIL: [brianbader@kpmg.com](mailto:brianbader@kpmg.com)

JUN YAN  
UNIVERSITY OF CONNECTICUT  
DEPARTMENT OF STATISTICS  
215 GLENBROOK RD. U-4120  
STORRS, CT 06269-4120  
E-MAIL: [jun.yan@uconn.edu](mailto:jun.yan@uconn.edu)

2 XUEBIN ZHANG  
ENVIRONMENT AND CLIMATE CHANGE CANADA  
CLIMATE RESEARCH DIVISION  
4905 DUFFERIN STREET  
DOWNSVIEW, ONTARIO M5H 5T4, CANADA  
E-MAIL: [xuebin.zhang@canada.ca](mailto:xuebin.zhang@canada.ca)

TABLE 1

*Empirical rejection rates of four goodness-of-fit tests for GPD under various data generation schemes with nominal size 0.05. GPDMix(a, b) refers to a 50/50 mixture of GPD(1, a) and GPD(1, b).*

Sample Size	50				100			
	Score	Moran	AD	CVM	Score	Moran	AD	CVM
Gamma(2, 1)	7.6	9.7	47.4	43.5	8.0	14.3	64.7	59.7
LogNormal	6.0	5.9	13.3	8.6	5.4	8.2	28.3	23.4
Weibull(0.75)	11.5	7.8	55.1	23.5	12.1	9.5	65.1	39.4
Weibull(1.25)	6.6	11.3	29.1	27.3	5.6	12.5	20.8	19.2
GPDMix(-0.4, 0.4)	11.4	7.5	19.2	9.9	16.0	8.6	24.3	20.4
GPDMix(0, 0.4)	7.8	6.0	6.5	5.9	7.4	5.9	9.6	5.6
GPDMix(-0.25, 0.25)	8.1	6.5	6.0	7.4	8.9	6.5	11.1	8.4
GPD(1, 0.25)	6.9	5.5	6.7	6.1	5.0	5.2	5.2	5.2
Sample Size	200				400			
Test	Score	Moran	AD	CVM	Score	Moran	AD	CVM
Gamma(2, 1)	15.4	23.3	95.3	93.1	36.5	42.2	100.0	100.0
LogNormal	5.7	11.9	69.3	59.7	7.9	19.1	97.8	95.0
Weibull(0.75)	16.5	10.7	84.8	66.4	32.1	14.6	98.2	93.0
Weibull(1.25)	8.0	14.7	40.9	36.7	15.5	19.2	79.8	74.6
GPDMix(-0.4, 0.4)	31.5	9.7	45.1	44.0	63.8	11.9	79.9	80.2
GPDMix(0, 0.4)	8.9	6.5	8.8	7.3	12.3	6.2	10.8	10.3
GPDMix(-0.25, 0.25)	13.9	6.7	16.6	14.8	26.2	7.9	33.0	32.4
GPD(1, 0.25)	5.3	5.8	7.2	5.2	4.7	5.3	5.8	4.7

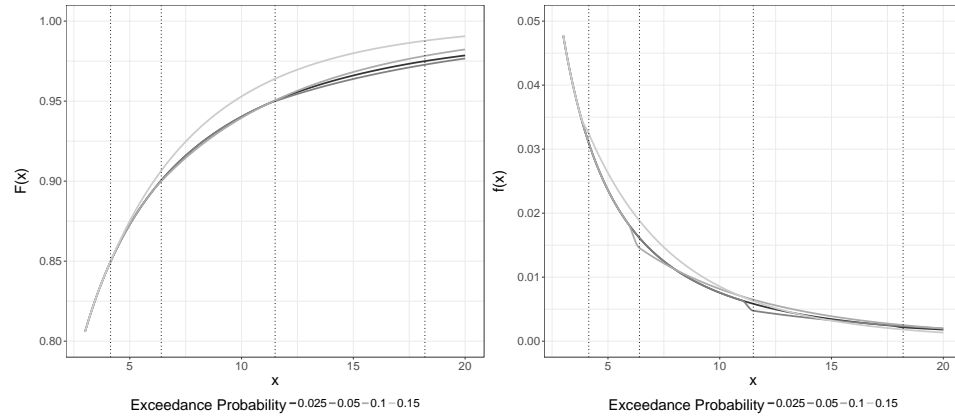


Fig 1: The cumulative and probability distributions used to generate data in the simulation study in Section 4. It is a Weibull distribution, with GPD tail. The start of the GP tail ( $u + \epsilon$ ) is 18.20, 11.50, 6.41, and 4.14 for the 0.025, 0.05, 0.10, and 0.15 exceedance probabilities, respectively.

1

2

3

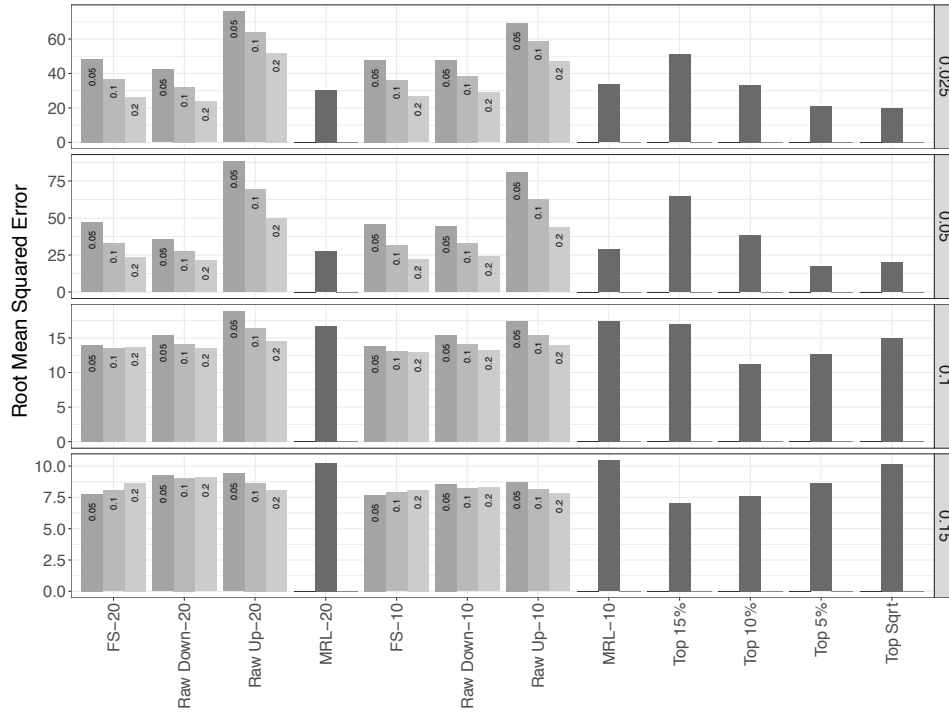


Fig 2: Root mean squared error (RMSE) of the estimators for the 50 year return level for the  $B = 1,500$  simulations and each method described in Section 4. The attached -10 and -20 refer to the two sets of thresholds tested. The decimals on the right axis reference the exceedance probability. The decimals in the bars of Raw and FS (ForwardStop) correspond to the significance level  $\alpha$  used in conjunction. Note that the ‘top’ methods are at a fixed threshold, thus no testing procedure is needed.

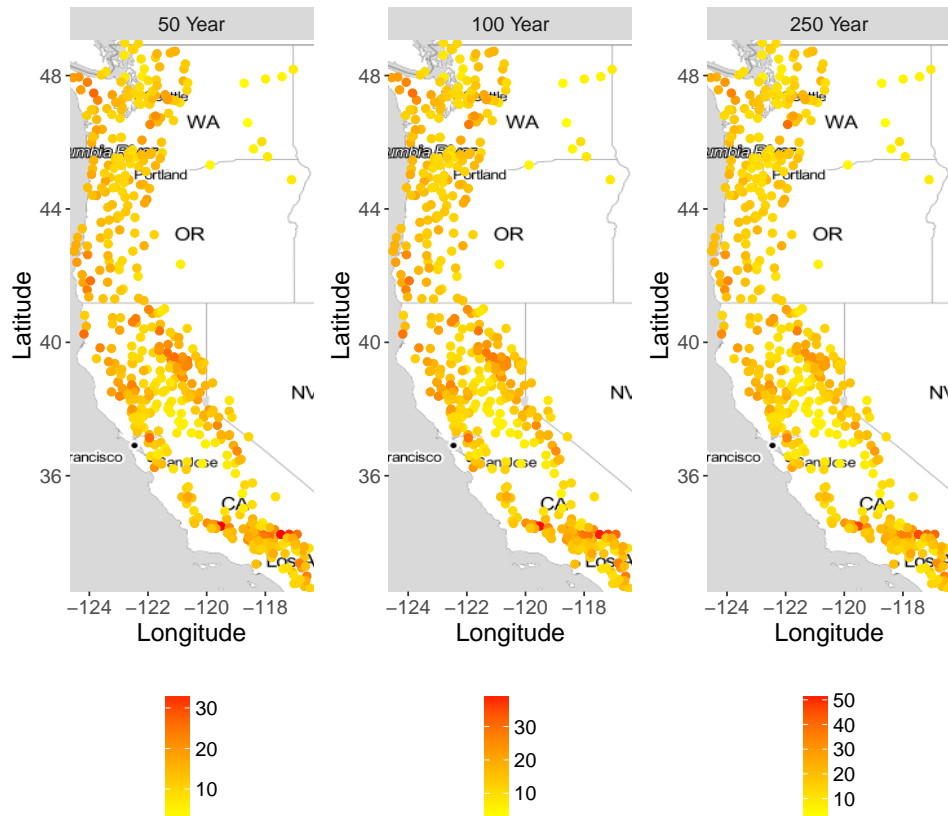


Fig 3: Estimated 50, 100, and 250 year daily precipitation return levels (cm) using the ForwardStop procedure and AD test, combined with jittering.