



HAL
open science

Extreme value modelling of SARS-CoV-2 community transmission using discrete Generalised Pareto distributions

Abdelaati Daouia, Gilles Stupfler, Antoine Usseglio-Carleve

► To cite this version:

Abdelaati Daouia, Gilles Stupfler, Antoine Usseglio-Carleve. Extreme value modelling of SARS-CoV-2 community transmission using discrete Generalised Pareto distributions. 2022. hal-03392044v2

HAL Id: hal-03392044

<https://hal.science/hal-03392044v2>

Preprint submitted on 27 Jul 2022 (v2), last revised 14 Mar 2023 (v4)

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Extreme value modelling of SARS-CoV-2 community transmission using discrete Generalised Pareto distributions

Abdelaati Daouia^a, Gilles Stupfler^b & Antoine Usseglio-Carleve^c

^a University of Toulouse Capitole, TSE - Decision Mathematics and Statistics, France
(ORCID: 0000-0003-2621-8860)

^b Univ Rennes, Ensai, CNRS, CREST - UMR 9194, F-35000 Rennes, France (ORCID:
0000-0003-2497-9412)

^c Avignon Université, Laboratoire de Mathématiques d'Avignon EA 2151, 84000 Avignon,
France (ORCID: 0000-0002-8148-3758)

Abstract. Superspreading has been suggested to be a major driver of overall transmission in the case of SARS-CoV-2. It is therefore important to statistically investigate the tail features of superspreading events (SSEs) to better understand virus propagation and control. Our extreme value analysis of different sources of secondary case data indicates that SSEs associated with SARS-CoV-2 may be fat-tailed, although substantially less so than predicted recently in the literature, but also less important relative to SSEs associated with SARS-CoV. The results caution against pooling data from both coronaviruses. This could provide policy- and decision-makers with a more reliable assessment of the tail exposure to SARS-CoV-2 contamination. Going further, we consider the broader problem of large community transmission. We study the tail behaviour of SARS-CoV-2 cluster cases documented both in official reports and in the media. Our results suggest that the observed cluster sizes have been fat-tailed in the vast majority of surveyed countries. We also give estimates and confidence intervals of the extreme potential risk for those countries. A key component of our methodology is up-to-date discrete Generalised Pareto models which allow for maximum-likelihood based inference of data with a high degree of discreteness.

Keywords. COVID-19, Superspreading, Cluster size, Secondary cases, Extreme value theory, Discrete extremes.

Introduction

Superspreading events (SSEs) have been recognised as a significant source of disease transmission for respiratory coronaviruses such as SARS-CoV and SARS-CoV-2 [1, 2]. SSEs may be defined as outbreaks in which a given individual (the index case)

24 infects a number of people (secondary cases) well above a certain measure, such
25 as the average or median number of infections. The number of secondary cases
26 resulting directly from an index case can be viewed as a random variable, say Z ,
27 defining the so-called offspring distribution. For both coronaviruses, events having
28 triggered more than 6 secondary cases have been suggested to constitute SSEs [3].
29 Data on such SSEs that was curated and reported in [3] in the early stages of the
30 COVID-19 pandemic is necessarily scarce: it consists mainly of 15 SSEs associated
31 with SARS-CoV and 45 SSEs associated with SARS-CoV-2, each represented by a
32 number of secondary cases Z_i resulting from a single given index case in Europe,
33 Asia or North America. The natural framework for the analysis of SSEs, and more
34 generally of atypical observations far away from the mean, is extreme value theory.
35 Following this framework, it was argued in [3] that SSEs are fat-tailed, although
36 this was done by pooling the 60 available SSEs from SARS-CoV and SARS-CoV-
37 2. A careful investigation of these SARS-CoV and SARS-CoV-2 datasets reveals
38 that the two largest observations in the pooled data are SARS-CoV SSEs; given the
39 small sample size, one may wonder whether the reported estimate of tail heaviness
40 is representative of the tail behaviour of SARS-CoV-2 SSEs.

41 This constitutes the motivation for this work, whose overarching goals are to
42 show how to conduct a principled extreme value analysis of community transmission
43 parameters, and to carry out such an analysis in the example of SARS-CoV-2. By
44 focusing directly on the raw SARS-CoV-2 data considered in [3], we provide evidence
45 of a lighter upper tail for SSEs with significantly less tail exposure than predicted in
46 their study. We arrive at the same conclusion by making use of a more recent and
47 much larger publicly available surveillance and contact-tracing database containing
48 the number of secondary cases Z_i for 88,527 index cases in the Indian states of
49 Andhra Pradesh and Tamil Nadu [4]. We also analyse two other South Korean
50 contact-tracing datasets, one collected in the first half of 2020 [3], the other during
51 the summer of 2021 when the Delta variant of SARS-CoV-2 was responsible for the
52 majority of positive cases [5]. The fat-tailedness of the secondary cases distribution
53 is found to be rather clear in the 2021 sample of data, while the analysis of the
54 2020 data is less conclusive. In all these samples of data we find point estimates
55 of the extreme value index suggesting that the secondary cases distribution has a
56 finite third moment, which stands in contrast with the earlier finding of [3] of a
57 distribution with an infinite variance.

58 In addition to that, we consider the broader problem of large community trans-
59 mission, as it represents the other fundamental source of pandemic risk. Large
60 infection clusters, along with SSEs, have been argued to play an important role in
61 the transmission of SARS-CoV-2 [2]. In a similar spirit to [2], we define a cluster
62 of SARS-CoV-2 cases in our analysis as a local outbreak involving a minimum of
63 two cases, including confirmed close contacts with epidemiological linkage over a
64 limited period of time. We consider two databases constructed from government
65 reports [6, 7, 8, 9] and media sources [10], comprising 15 samples of SARS-CoV-2
66 cluster sizes recorded in 11 countries and 4 US states. Our results show that 13
67 of these 15 countries and states have fat-tailed cluster size distributions, thus fa-
68 cilitating the process of inferring their risk category in terms of large community
69 transmission. This allows us to better understand the drivers of superspreading and

70 cluster formation in the ongoing COVID-19 pandemic. The recent theory of discrete
71 extremes [11, 12, 13, 14] is our basic tool to address the highly discrete nature of
72 SARS-CoV-2 secondary transmission data and cluster sizes. Its use constitutes our
73 main statistical contribution to the study of the transmission of the SARS-CoV-2
74 virus. As we illustrate throughout the paper, estimating and inferring the extreme
75 value index and extreme percentiles of the underlying discrete distributions with this
76 methodology is much easier and accurate than with classical extreme value meth-
77 ods such as the Hill and Generalised Pareto maximum likelihood estimators, which
78 heavily rely on the continuous data assumption.

79 The structure of the paper is as follows. We first describe the methods employed
80 throughout our study, including the discrete Generalised Pareto Distribution fitted
81 to exceedances over a high threshold by means of the maximum likelihood estima-
82 tor. We then analyse our datasets, first on SARS-CoV-2 secondary case numbers
83 and then on cluster sizes, using these methods. A Discussion section gathers and
84 contrasts these findings and concludes with additional comments about the scope,
85 limitations and robustness of our results, as well as ideas for further work.

86 Methods

87 We use several methods from extreme value theory, which constitutes the correct
88 mathematical framework for the analysis of high observations from a random phe-
89 nomenon [15]. We are particularly interested in methods that can describe so-called
90 heavy-tailed random variables, which infrequently but regularly generate very high
91 values and therefore appear to be relevant in the analysis of SARS-CoV-2 transmis-
92 sion. A random variable X is heavy-tailed (or fat-tailed) if and only if its distribution
93 function $\mathbb{P}(X \leq x)$ can be expressed as $\mathbb{P}(X \leq x) = 1 - x^{-1/\xi}\ell(x)$, where ℓ sat-
94 isfies $\ell(tx)/\ell(t) \rightarrow 1$ as $t \rightarrow \infty$ for any positive real number x . Informally, the
95 tail behaviour of X is controlled by the extreme value index $\xi > 0$, which must be
96 estimated to get a precise understanding of tail heaviness. A standard estimator in
97 this context is the Hill estimator [16]. For a dataset Z_1, \dots, Z_n , the Hill estimator
98 at threshold u is defined as

$$\hat{\xi}_u^H = \frac{1}{\sum_{i=1}^n \mathbb{1}_{\{Z_i > u\}}} \sum_{i=1}^n \log \left(\frac{Z_i}{u} \right) \mathbb{1}_{\{Z_i > u\}}.$$

99 It is of course crucial, before using the Hill estimator, to ascertain whether the
100 distribution of the data points indeed has a heavy tail. A common diagnostic method
101 is the mean excess plot, which estimates the values of the mean excess function
102 $E(u) = \mathbb{E}[Z - u | Z > u]$ as function of u . A natural estimate of $E(u)$ is given, for
103 each threshold u , by its empirical counterpart

$$\hat{E}(u) = \frac{\sum_{i=1}^n Z_i \mathbb{1}_{\{Z_i > u\}}}{\sum_{i=1}^n \mathbb{1}_{\{Z_i > u\}}} - u.$$

A heavy-tailed distribution will typically have mean excess plots exhibiting a linear
upward drift for large values of u , see for example Section 1.2.2 in [17]. It has, how-
ever, been observed in the extreme value literature [18] that the mean excess function

very often exhibits a non-linear behaviour at the right end of the mean excess plot, due to very high variability of the estimate of $E(u)$ when u is close to the highest Z_i . As a consequence, good statistical practice recommends to confirm a diagnostic of a heavy tail using other extreme value tools. One such general approach, which does not presuppose that the data is heavy-tailed, consists in using the Generalised Pareto maximum likelihood estimator, defined as, according to Section 5.3.2 in [17]:

$$\begin{aligned} & \left(\hat{\xi}_u^{GP}, \hat{\sigma}_u^{GP} \right) \\ &= \arg \min_{(\xi, \sigma) \in (-1/2, \infty) \times (0, \infty)} \sum_{i=1}^n \left[-\log \sigma - \left(\frac{1}{\gamma} + 1 \right) \log \left(1 + \xi \frac{Z_i - u}{\sigma} \right) \right] \mathbb{1}_{\{Z_i > u\}}. \end{aligned}$$

104 The Generalised Pareto maximum likelihood estimators are valid even when the
105 underlying distribution is not heavy-tailed, which has made them very popular in
106 the natural sciences [19].

107 However, both the Hill and Generalised Pareto estimators of ξ suffer from jagged
108 sample paths when the data points Z_i come from a distribution with a high degree
109 of discreteness. This behavior makes it extremely difficult to choose an accurate
110 estimate of ξ , which renders the two methods highly unsatisfactory. The essential
111 reason behind this phenomenon is that both estimators are built under the – gener-
112 ally incorrect – assumption that the data points come from a pure (Generalised)
113 Pareto distribution, which is continuous, and as such, they cannot be expected to
114 handle a substantial degree of discreteness. We exemplify this phenomenon in Fig. 1:
115 notice, in the top panels, the stark difference in stability and smoothness of sam-
116 ple paths between a Hill plot for continuous data Z_i and its counterpart for data
117 rounded to the nearest integer up. The bottom panels show that the Hill estima-
118 tor for discrete data tends to be strongly biased and much more so than the Hill
119 estimator for continuous data.

A statistically principled alternative is to employ proper discrete models to con-
struct an estimator of the extreme value index. This was pursued by [13], which used
so-called D-GPD (for Discrete-Generalised Pareto Distribution) models to introduce
the maximum likelihood estimators

$$\begin{aligned} & \left(\hat{\xi}_u, \hat{\sigma}_u \right) \\ &= \arg \min_{(\xi, \sigma) \in \mathbb{R} \times (0, \infty)} \sum_{i=1}^n \log \left(\left(1 + \xi \frac{Z_i - u}{\sigma} \right)^{-1/\xi} - \left(1 + \xi \frac{Z_i - u + 1}{\sigma} \right)^{-1/\xi} \right) \mathbb{1}_{\{Z_i \geq u\}}. \end{aligned}$$

120 When $\xi = 0$, the convention we adopt is that $(1 + \xi z)^{-1/\xi} = \exp(-z)$, for any
121 $z \in \mathbb{R}$. These maximum likelihood estimators of the extreme value index ξ and scale
122 parameter σ of the D-GPD model are readily obtained through the R maximisa-
123 tion routine `optim`. Using the classical theory of maximum likelihood estimators,
124 confidence intervals for ξ may be derived from $\hat{\xi}_u$ by estimating the total Fisher
125 information matrix $I(\xi, \sigma)$ using a finite difference method and then deducing the
126 following $100\alpha\%$ -confidence interval for ξ :

$$\left[\hat{\xi}_u + \sqrt{\left(\hat{I}(\xi, \sigma)^{-1} \right)_{1,1}} \Phi^{-1} \left(\frac{1 - \alpha}{2} \right), \hat{\xi}_u + \sqrt{\left(\hat{I}(\xi, \sigma)^{-1} \right)_{1,1}} \Phi^{-1} \left(\frac{1 + \alpha}{2} \right) \right],$$

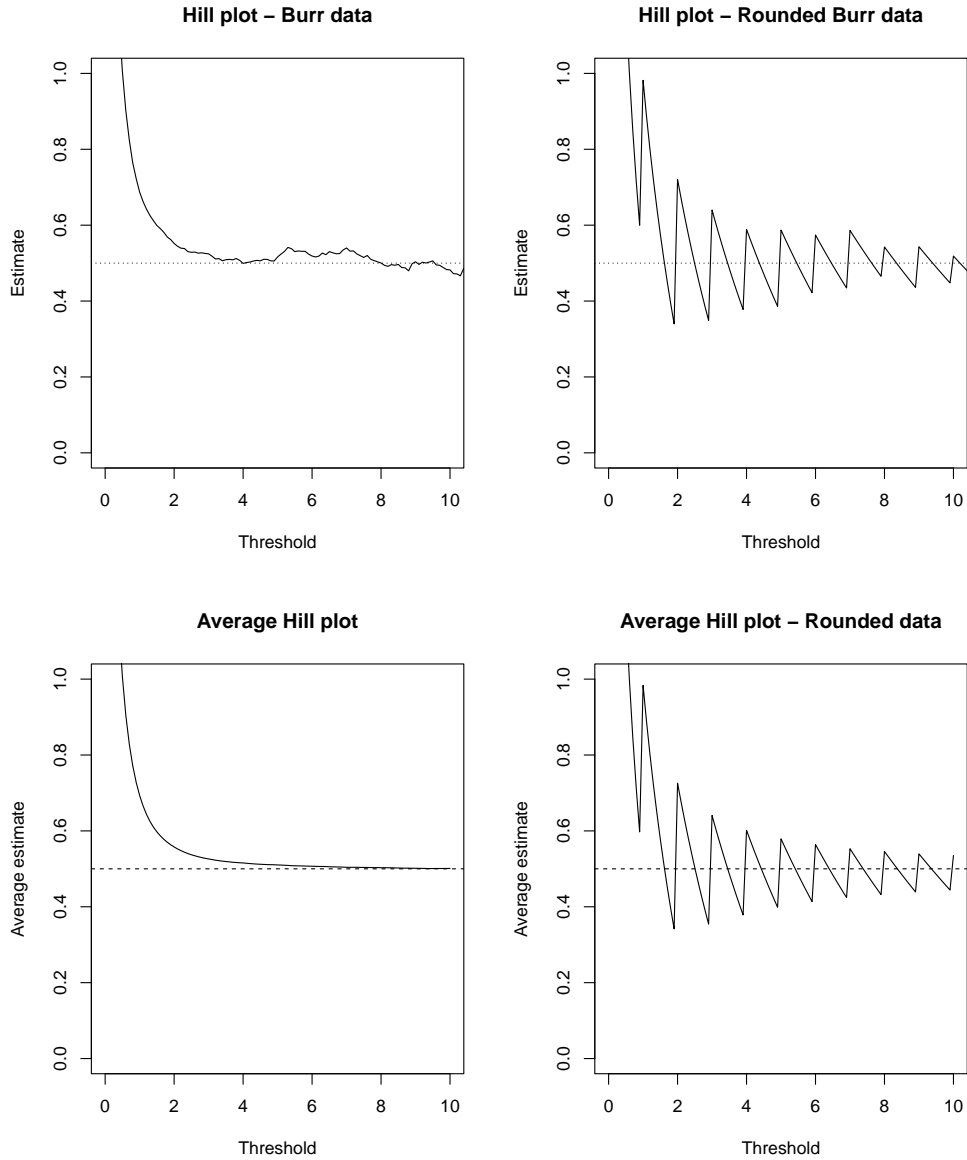


Figure 1: Top panels: Hill plots as functions of the threshold value u , for $n = 10,000$ simulated data points Z_i from the Burr distribution with probability density function $f(x) = \xi^{-1}x^{-\rho/\xi-1}(1+x^{-\rho/\xi})^{1/\rho-1}$ (for $x > 0$) with $\xi = 1/2$ and $\rho = -1$ in the left panel, and for the data $\lceil Z_i \rceil$ (*i.e.* the smallest integer larger than or equal to Z_i) in the right panel. Bottom panels: Averaged Hill plots when this experiment is repeated $N = 1,000$ times.

127 where Φ denotes the standard normal distribution function and Φ^{-1} its inverse
 128 (quantile function). Modelling $Z - u$ conditional on $Z \geq u$ by a D-GPD distribution
 129 with parameter estimates $(\hat{\xi}_u, \hat{\sigma}_u)$ suggests the following estimate of the 100 α th
 130 percentile of Z adapted from [12, Formula (5) p.41]:

$$\hat{q}_\alpha = \left\lceil \frac{\hat{\sigma}_u}{\hat{\xi}_u} \left(\left(\frac{n(1-\alpha)}{\sum_{i=1}^n \mathbb{1}_{\{Z_i \geq u\}}} \right)^{-\hat{\xi}_u} - 1 \right) + u - 1 \right\rceil,$$

131 for α large enough. Here, $\lceil \cdot \rceil$ denotes the ceiling function, that is, $\lceil x \rceil$ denotes the
 132 smallest integer larger than or equal to x . Estimating this quantile by plugging in
 133 the aforementioned estimates of ξ and σ makes it possible to infer extreme quantile
 134 levels and therefore get precise information on the tail behaviour of a distribution
 135 with a large degree of discreteness.

136 For comparison purposes, we will contrast the resulting extreme quantile estimates
 137 with those provided by the (conditioned) negative binomial distribution. Recall that
 138 the probability mass function of the negative binomial distribution (with parameters
 139 $r > 0$ and $p \in (0, 1)$) conditional on $Z > u$, is given by

$$\mathbb{P}_{p,r,u}(Z = k) = \frac{\frac{\Gamma(k+r)}{k! \Gamma(r)} p^r (1-p)^k}{1 - \sum_{i=0}^u \frac{\Gamma(i+r)}{i! \Gamma(r)} p^r (1-p)^i}, \text{ for all } k > u.$$

140 Here Γ denotes Euler's Gamma function. With a dataset z_1, \dots, z_n , the parameter
 141 estimators are therefore obtained as the maximum log-likelihood solution

$$\arg \max_{(p,r) \in (0,1) \times (0,\infty)} \sum_{i=1}^n \log \mathbb{P}_{p,r,u}(Z = z_i).$$

142 Ever since the seminal work of [1], the negative binomial distribution has been widely
 143 used to describe the number of secondary cases resulting from an index case of SARS-
 144 CoV. As suggested in [3, 21], this model has exponentially decreasing probability
 145 mass functions and thus cannot be expected to accurately represent tail heaviness in
 146 SARS-CoV-2 transmission data. We provide below further evidence for this claim,
 147 and for the suitability of D-GPD maximum likelihood estimates in the context of
 148 discrete data, through several datasets gathering numbers of SARS-CoV-2 secondary
 149 cases and cluster sizes in different settings.

150 Data and results

151 **Analysis of secondary case data.** Our first two datasets were reported in [3].
 152 They consist of 15 SSEs associated with SARS-CoV (Dataset S1) and 45 SSEs
 153 associated with SARS-CoV-2 (Dataset S2), each resulting in more than 6 secondary
 154 cases, along with month of occurrence and location of the superspreading event,
 155 and its setting. We refer to [3] for further details about the construction of these
 156 datasets. Pooling the 15 SSEs associated with SARS-CoV and 45 SSEs associated
 157 with SARS-CoV-2 into a single sample and making use of a Generalised Pareto
 158 approximation, [3] has suggested that the distribution of the number of secondary

159 cases Z belongs to the Fréchet maximum domain of attraction [20], that is, the
 160 set of Pareto-type distributions, with tail index ξ between 0.5 and 1 (the estimate
 161 provided in [3, Fig. 1 E] is $\hat{\xi} \approx 0.6$). The index ξ tunes the tail heaviness of the
 162 distribution, with higher positive values indicating a heavier upper tail: moments of
 163 order higher than or equal to $1/\xi$ do not exist. An estimate of ξ around 0.6 means
 164 that the second moment of Z does not exist, reflecting the outsized contribution
 165 of SSEs to overall transmission. Most importantly perhaps, these findings on the
 166 tail heaviness of Z invalidate the conventional assumption that Z follows a negative
 167 binomial distribution for either coronavirus, whereas this assumption was widely
 168 adopted in the literature on disease transmission ever since the influential work [1]
 169 on SARS-CoV, and it is still widely employed for SARS-CoV-2, see *e.g.* [5, 22, 23].

170 Based on our statistical analysis of these datasets, summarised in Fig. 2, one
 171 may however argue that the method of [3] is inappropriate for examining the tail
 172 behaviour of their particular 60 SSEs. The sparsity of data on SSEs is addressed by
 173 combining the 15 and 45 observations associated with SARS-CoV and SARS-CoV-2
 174 into a single sample, whereas the two datasets correspond to completely different
 175 distributions (Fig. 2 (a)) and should not be pooled accordingly. This is apparent
 176 from either a Kolmogorov-Smirnov test, with p -value 0.015, or the more common
 177 approach making the questionable assumption that Z follows a negative binomial
 178 distribution. The conditional (given $Z > 6$) negative binomial fit of the probability
 179 mass function to the Z_i (by construction larger than 6), calculated as described
 180 in the last paragraph of the Methods section (Fig. 2 (b)), already suggests that
 181 the upper tail of Z for SARS-CoV appreciably dominates that for SARS-CoV-2.
 182 In other words, even a naive analysis of the SSE distributions, using the classical
 183 negative binomial distribution and not accounting for the heavy tail in the data,
 184 indicates that the SSEs for SARS-CoV and those for SARS-CoV-2 exhibit different
 185 statistical behaviour. This is confirmed by a proper extreme value analysis of the
 186 data (Fig. 2 (c)): the ξ estimates obtained from the Hill estimator in the special
 187 case of SARS-CoV-2 vary between 0.35 and 0.45, and as such differ substantially
 188 from the various competing estimates found to vary between 0.5 and 1 in [3]. Even
 189 the 90% confidence intervals of ξ for SARS-CoV-2 (dashed red lines in Fig. 2 (c))
 190 only partially contain the estimated tail index plot for SARS-CoV (solid blue line),
 191 reflecting a net difference between the two heavy-tailed distributions of secondary
 192 cases associated with SARS-CoV and SARS-CoV-2. This conclusion is corroborated
 193 by the mean excess function estimates (Fig. 2 (d)), which similarly indicate the
 194 relevance of separating the analysis for each coronavirus. This suggests that although
 195 SARS-CoV and SARS-CoV-2 belong to the same family of respiratory diseases,
 196 superspreading events are larger in scale for SARS-CoV in comparison to SARS-
 197 CoV-2. For all these reasons, pooling the data before applying extreme value tools
 198 can lead to misleading conclusions on the propagation of the SARS-CoV-2 virus.

199 Yet, the low sample size of this SSE dataset puts a question mark over the quality
 200 of the statistical analysis. Trustworthy extreme value inference may require a larger
 201 sample size, of the order of at least several thousands. This is why we also analysed
 202 a much larger Indian secondary case dataset of size $n = 88,527$ (Database S3). This
 203 comprehensive surveillance and contact-tracing database was collected in 2020 by
 204 the public health authorities of the two Indian states of Andhra Pradesh and Tamil

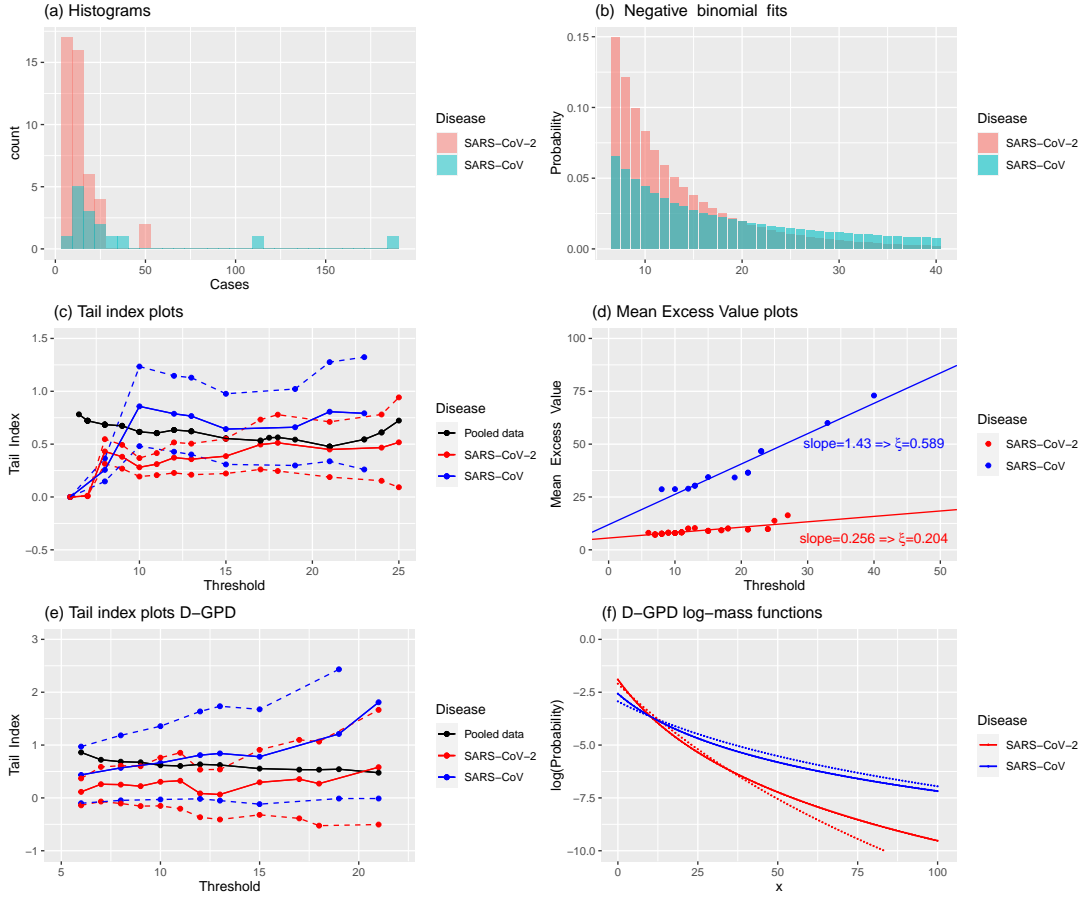


Figure 2: Secondary case data from [3] (Datasets S1 and S2). (a) Histogram of the number of secondary cases for SARS-CoV (blue, $n = 15$) and SARS-CoV-2 (red, $n = 45$) SSEs. (b) Fitted probability mass function, conditional on $Z > 6$, of the negative binomial distribution for SARS-CoV (blue) and SARS-CoV-2 (red) SSEs. (c) Hill estimates of ξ for SSEs associated with SARS-CoV (solid blue), SARS-CoV-2 (solid red), and the pooled data (solid black), obtained from the exceedance values $Z_i - u$ given $Z_i \geq u$, as function of the threshold u , along with the resulting 90% confidence intervals for SARS-CoV (dashed blue) and SARS-CoV-2 (dashed red) SSEs. (d) Mean excess plots of SARS-CoV (blue) and SARS-CoV-2 (red) SSEs, quantified by the average of the exceedances $Z_i - u$ given $Z_i \geq u$, as function of u . (e) Discrete GPD maximum likelihood estimates of ξ for SARS-CoV (solid blue) and SARS-CoV-2 (solid red) SSEs, calculated from the exceedances $Z_i - u$ given $Z_i \geq u$, as function of u , along with their corresponding 90% confidence intervals (dashed lines), and the Hill plot produced by combining SARS-CoV and SARS-CoV-2 SSEs. (f) Logarithm of the probability mass functions $\mathbb{P}_{\sigma,\xi}(X = x)$ of the D-GPD fits to the exceedance values $Z_i - u$ given $Z_i \geq u$, for the thresholds $u = 6$ (dotted lines) and $u = 10$ (solid lines), for SARS-CoV (blue) and SARS-CoV-2 (red).

205 Nadu, whose residents total about 10% of India’s population. It was studied for
 206 instance in [4] and [21], and we refer to the latter for more information about the
 207 database’s construction and contents. Results are reported in Fig. 3. Although the
 208 barplot of this data (Fig. 3 (a)) gives evidence of a considerable right skewness and
 209 its summary extreme value analysis (Fig. 3 (b)) suggests a heavy right tail, it should
 210 be noted that since the Z_i range from 0 to 39 with a sample size of 88,527, the data
 211 is necessarily highly discrete with a large number of tied observations (see Table 1).

Z	0	1	2	3	4	5	6	7	8	9	10	11				
Count	62,540	17,493	4,885	1,730	802	444	267	149	67	44	29	22				
Z	12	13	14	15	16	17	18	19	21	22	23	25	28	31	37	39
Count	14	16	3	3	4	4	1	1	2	1	1	1	1	1	1	1

Table 1: Secondary case data (Database S3) for SARS-CoV-2 from Andhra Pradesh and Tamil Nadu (India).

212 Ignoring the discrete nature of the Z_i by modelling their tail behaviour with the
 213 (Generalised) Pareto distribution is inappropriate as this typically results in unreli-
 214 able tail index estimates and confidence intervals [13]. This becomes obvious here by
 215 superimposing both the classical Hill and continuous Generalised Pareto maximum
 216 likelihood estimators of the extreme value index, as functions of a varying thresh-
 217 old u in Fig. 3 (c). Clearly, both plots are so volatile and jagged that it is hard
 218 to identify any stable region and therefore a reasonable point estimate of ξ cannot
 219 easily be determined. We address this limitation by applying the recent theory of
 220 discrete extremes developed in [11, 13] and based on the discrete Generalised Pareto
 221 distribution (D-GPD). The D-GPD, first employed by [12] to model road accidents
 222 and more recently in [14] to model hospital congestion, has been shown to outper-
 223 form the continuous GPD when there are a large number of tied observations: see
 224 the simulated Poisson and discrete Inverse-Gamma examples in Section 3.1 of [13],
 225 which respectively show that the GPD provides poor fits and poor tail estimates
 226 when the data is highly discrete, while the D-GPD distribution performs well. Its
 227 closed-form survival and probability mass functions allow for an exact likelihood-
 228 based inference. Using the D-GPD distribution to fit exceedances $Z_i - u$ above the
 229 threshold u (rather than trying to fit the whole of the distribution, as [21] did using
 230 a discrete Pareto distribution) results in a much smoother and stable fit (Fig. 3 (c)),
 231 and leads to an estimate of ξ around 0.24 with the 90% confidence intervals over-
 232 whelmingly suggesting an estimate greater than 0, thus confirming the heavy-tailed
 233 nature of SARS-CoV-2 SSEs (Fig. 3 (d)) in this sample. Interestingly, revisiting the
 234 small SARS-CoV-2 SSE dataset (Dataset S2) of size 45 using the D-GPD maximum
 235 likelihood estimation method (Fig. 2 (e)) results in an estimate of around 0.25, in
 236 agreement with the results from the Indian secondary case data. This suggests that
 237 the distribution of SARS-CoV-2 SSEs has a finite third moment and possibly even
 238 a fourth moment. These results are different from those obtained for the SARS-
 239 CoV SSEs. The latter rather point towards a distribution with infinite variance and
 240 thus a much heavier right tail. This is confirmed by considering the fitted D-GPD
 241 probability mass functions for secondary cases (Fig. 2 (f)) that decrease much more
 242 rapidly for SARS-CoV-2 than for SARS-CoV.

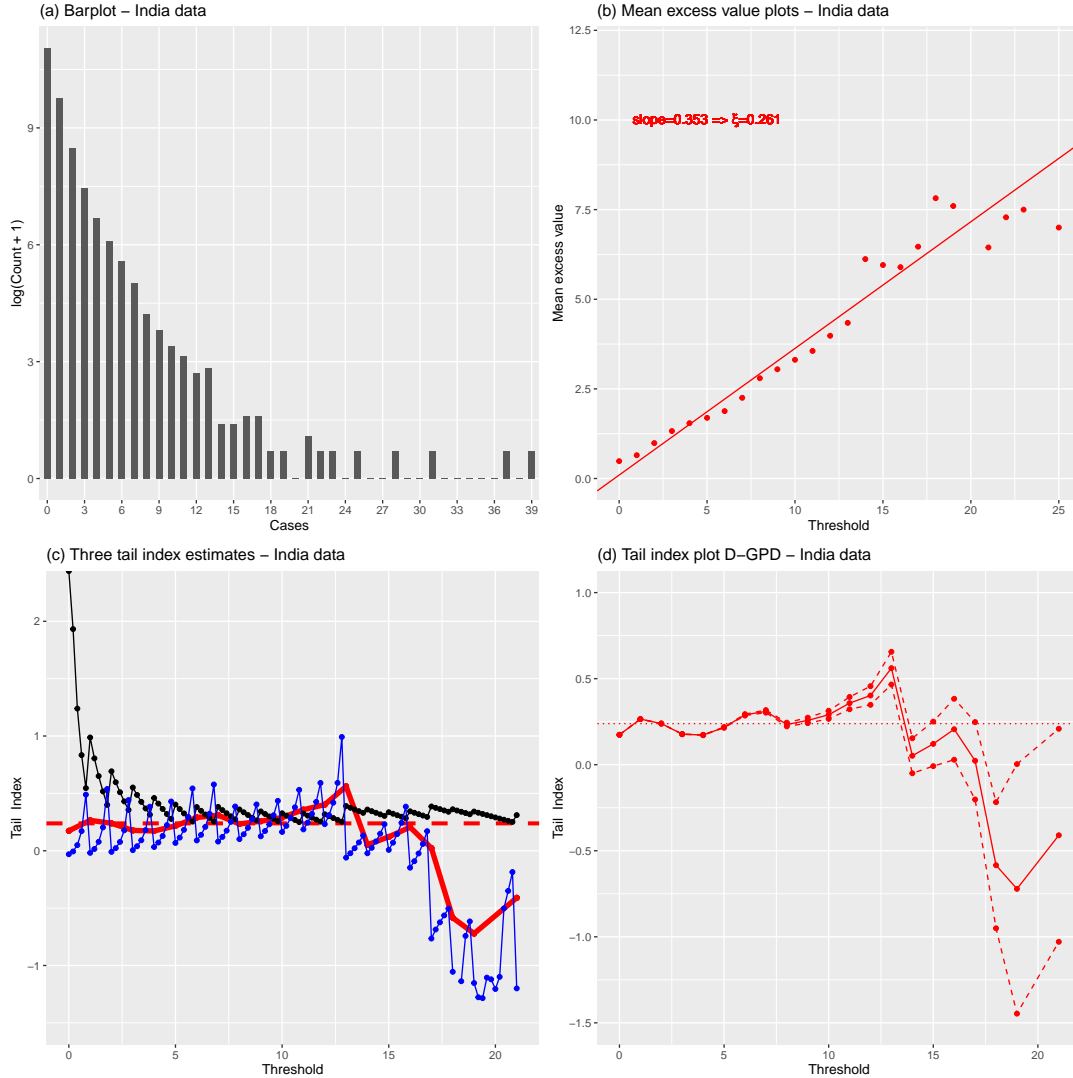


Figure 3: Secondary case data (Database S3) for SARS-CoV-2 from Andhra Pradesh and Tamil Nadu (India). (a) Bar plot of the $\log(Z_i+1)$ ($n = 88,527$). (b) Mean excess plots of secondary cases. (c) Hill (solid black), continuous GPD maximum likelihood (solid blue) and discrete GPD maximum likelihood (solid bold red) estimates of ξ . (d) Discrete GPD maximum likelihood estimates of ξ (solid red) and their associated 90% confidence intervals (dashed red). In panels (c) and (d), the averaged discrete GPD estimate $\hat{\xi} = 0.239$ over the stable region $u \in [0, 10]$ is indicated with the horizontal red line.

243 To examine the extreme value behaviour of the SARS-CoV-2 offspring distribu-
 244 tion in different conditions, we turn to the analysis of two contact-tracing datasets in
 245 South Korea, a country which has a similar population density to the Indian state of
 246 Tamil Nadu, but did not resort to any full lockdown and has one of the largest and
 247 best-organised epidemic control programmes in the world. The first dataset was col-
 248 lected in the first half of 2020 (Database S4), while the second was collected during
 249 the fourth community epidemic in the summer of 2021 (Database S5) in the context
 250 of the assessment of transmission dynamics for the Delta variant of SARS-CoV-2.
 251 The first dataset, which consists of $n = 5,165$ numbers of SARS-CoV-2 secondary
 252 cases Z_i , was analysed in [3]. See Table 2.

Z	0	1	2	3	4	5	6	7	8	9
Count	4,558	364	114	62	27	7	7	4	4	1
Z	10	11	12	15	17	18	21	24	27	51
Count	2	3	1	2	2	1	2	2	1	1

Table 2: Secondary case data (Database S4) for SARS-CoV-2 collected in South Korea in the first half of 2020.

253 We revisit the estimation of, and inference about, the underlying extreme value
 254 index by comparing the D-GPD estimates with the classical GPD and Hill estimates.
 255 Results are displayed in Fig. 4. A least squares fit to the first part of the mean excess
 256 plot (Fig. 4 (b)) suggests a linearly increasing fit to the mean excess function with
 257 a slope of around 0.85, but this ignores the flat or even slightly linearly decreasing
 258 right-hand part of the data cloud. This throws the assumption that the offspring
 259 distribution is heavy-tailed in doubt, although the barplot of the data (Fig. 4 (a))
 260 would tentatively back the heavy tail assumption. The Hill estimator, which pre-
 261 supposes that the data is heavy-tailed and graphed as a black line in Fig. 4 (c), does
 262 not exhibit any stable region which would allow to produce a reasonable point esti-
 263 mate. In such scenarios, best practice in extreme value theory requires calculating
 264 alternative extreme value estimators whose consistency does not rest upon the heavy
 265 tail assumption (unlike the Hill estimator), such as the general GPD and D-GPD
 266 estimators. These are also represented in Fig. 4 (c). Clearly, the paths of these two
 267 estimates follow a similar trajectory which is very different from that of the Hill
 268 plot. They point towards substantially lower estimates of ξ , and even though the
 269 estimates are overall larger than 0, the validity of the heavy tail assumption $\xi > 0$
 270 is not obvious for this dataset. Fig. 4 (d) further supports this observation: in the
 271 (somewhat) stable region around the threshold $u = 10$, the 90% confidence interval
 272 produced through maximum likelihood theory contains the value 0. Our conclusion
 273 from the analysis of this dataset is that the distribution of the number of secondary
 274 cases is either fat-tailed but with a low tail index, or perhaps even light-tailed. As
 275 a consequence, our finding is qualitatively different from that of [3], since we do not
 276 obtain ξ estimates similar to those found by merging Datasets S1 and S2.

277 The second South Korean contact-tracing dataset comprises $n = 33,903$ SARS-
 278 CoV-2 numbers of secondary cases Z_i (Database S5) detected between 25th July
 279 2021 and 15th August 2021. It was initially explored in [5], where it was highlighted
 280 that the Delta variant accounted for the majority of those cases. We therefore inves-

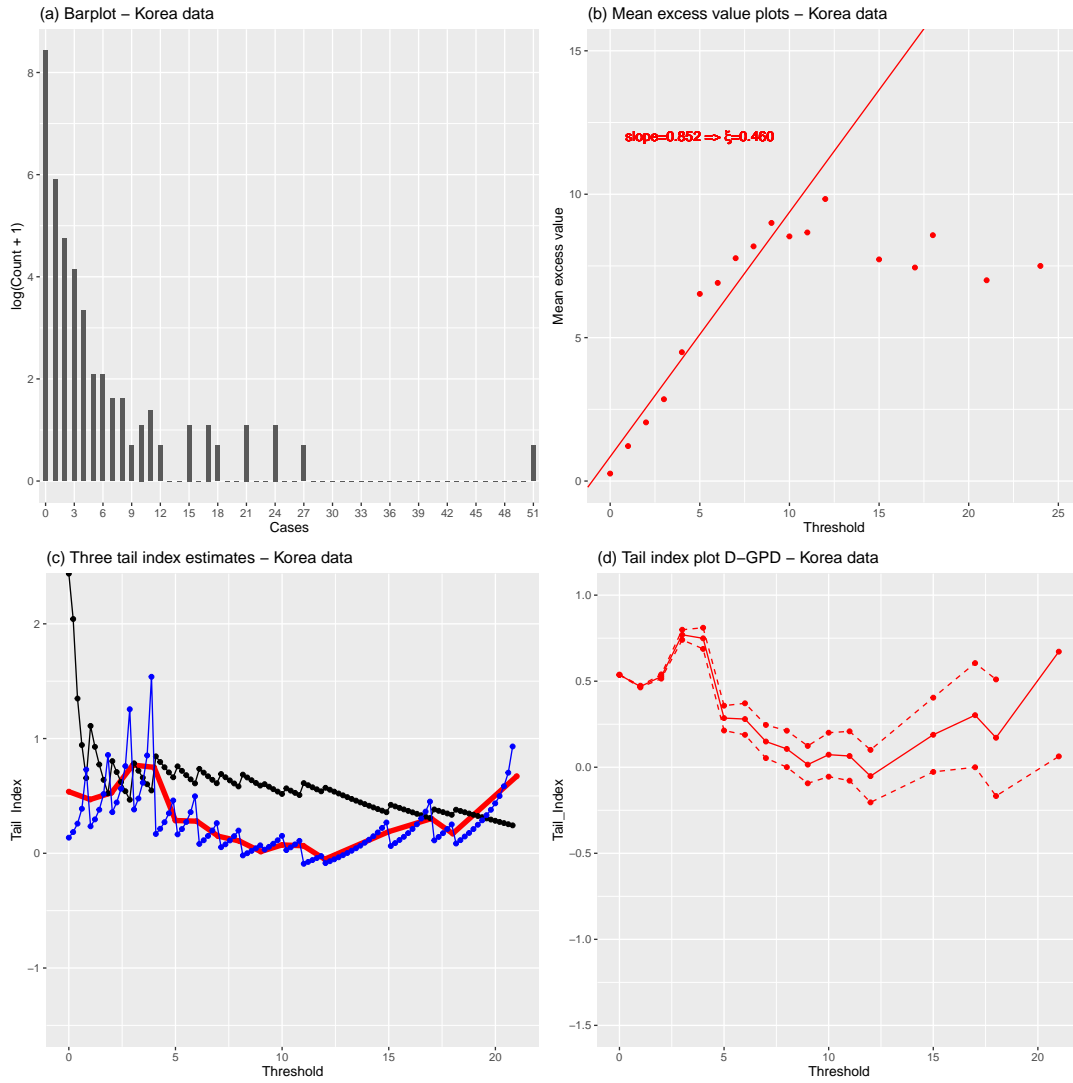


Figure 4: Secondary case data (Database S4) for SARS-CoV-2 from South Korea (first half of 2020). (a) Bar plot of the $\log(Z_i + 1)$ ($n = 5,165$). (b) Mean excess plots of secondary cases. (c) Hill (solid black), continuous GPD maximum likelihood (solid blue) and discrete GPD maximum likelihood (solid bold red) estimates of ξ . (d) Discrete GPD maximum likelihood estimates of ξ (solid red) and their associated 90% confidence intervals (dashed red).

281 tigate this dataset to ascertain whether the tail behaviour of SSEs is substantially
282 different for the Delta variant. The data is presented in Table 3 below. The re-
283 sults we obtain for this dataset are displayed in Fig. 5. The barplot of the data in
284 Fig. 5 (a) again backs the assumption of a heavy tail, but here, the mean excess plot
285 in Fig. 5 (b) suggests a more convincing linearly increasing fit to the mean excess
286 function with a slope of around 0.3. The Hill estimator and both continuous and
287 discrete GPD maximum likelihood estimators, represented in Fig. 5 (c), appear to
288 support the fat tail assumption of the offspring distribution which is mainly dom-
289 inated here by the Delta variant. Once again, the D-GPD estimate has a much
290 smoother and more stable sample path, with a stable zone over $u \in [1, 10]$ indicat-
291 ing a point estimate of around 0.21. The 90% confidence interval of the D-GPD
292 estimate over that region, provided in Fig. 5 (d), does not contain 0 and offers fur-
293 ther justification of the assumption that the offspring distribution is heavy-tailed
294 in this dataset, in contrast to the 2020 South Korea data where the validity of this
295 conclusion is much less clear.

Z	0	1	2	3	4	5	6	7	8	9	10		
Count	29,193	2,154	1,121	594	332	207	113	53	53	21	21		
Z	11	12	13	14	15	16	17	18	19	21	22	24	32
Count	6	8	5	3	3	2	2	3	3	1	2	2	1

Table 3: Secondary case data (Database S5) for SARS-CoV-2 collected in South Korea from 25th July 2021 to 15th August 2021.

296 **Analysis of cluster size data.** We broaden our analysis by examining whether
297 SARS-CoV-2 cluster sizes are fat-tailed. We consider a database of 15 samples of
298 cluster sizes recorded in 11 countries and 4 US states. We define a cluster as a
299 local outbreak involving a minimum of two cases, including confirmed close contacts
300 with epidemiological linkage over a limited period of time. The number of reported
301 clusters per country or state varies from 29 (France) to 4,769 (Colorado, USA). The
302 database is constructed from government reports [6, 7, 8, 9] (Database S6) and media
303 sources [10] (Database S7). The median cluster sizes were 5 (Database S6) and 33
304 (Database S7), and the largest clusters had sizes 1,761 (Database S6, in a Colorado
305 prison) and 7,000 (Database S7, in an Italian football stadium). We denote by Y_i
306 the number of SARS-CoV-2 cases in cluster i . The ξ estimates from each sample of
307 cluster sizes allow to infer the risk category of the corresponding country/state in
308 terms of local community transmission.

309 Figs. 6 and 7 display the D-GPD maximum likelihood estimates of ξ as functions
310 of the cluster size u . A common practice for selecting a suitable pointwise estimate of
311 ξ is to pick out a sufficiently high threshold u corresponding to a stable region of the
312 plot [15], as indicated by the vertical dashed lines in Figs. 6 and 7. The final selected
313 estimates are reported in Table 4, where 13 out of the 15 countries or states appear
314 to have fat-tailed cluster size distributions (confirmed at the 90% confidence level
315 except for China). The analysis for California and UK & Ireland was inconclusive.
316 For the California dataset, this is possibly due to a strong degree of heterogeneity
317 (see the histogram in the bottom left panel of Fig. 7). A stratified study of the

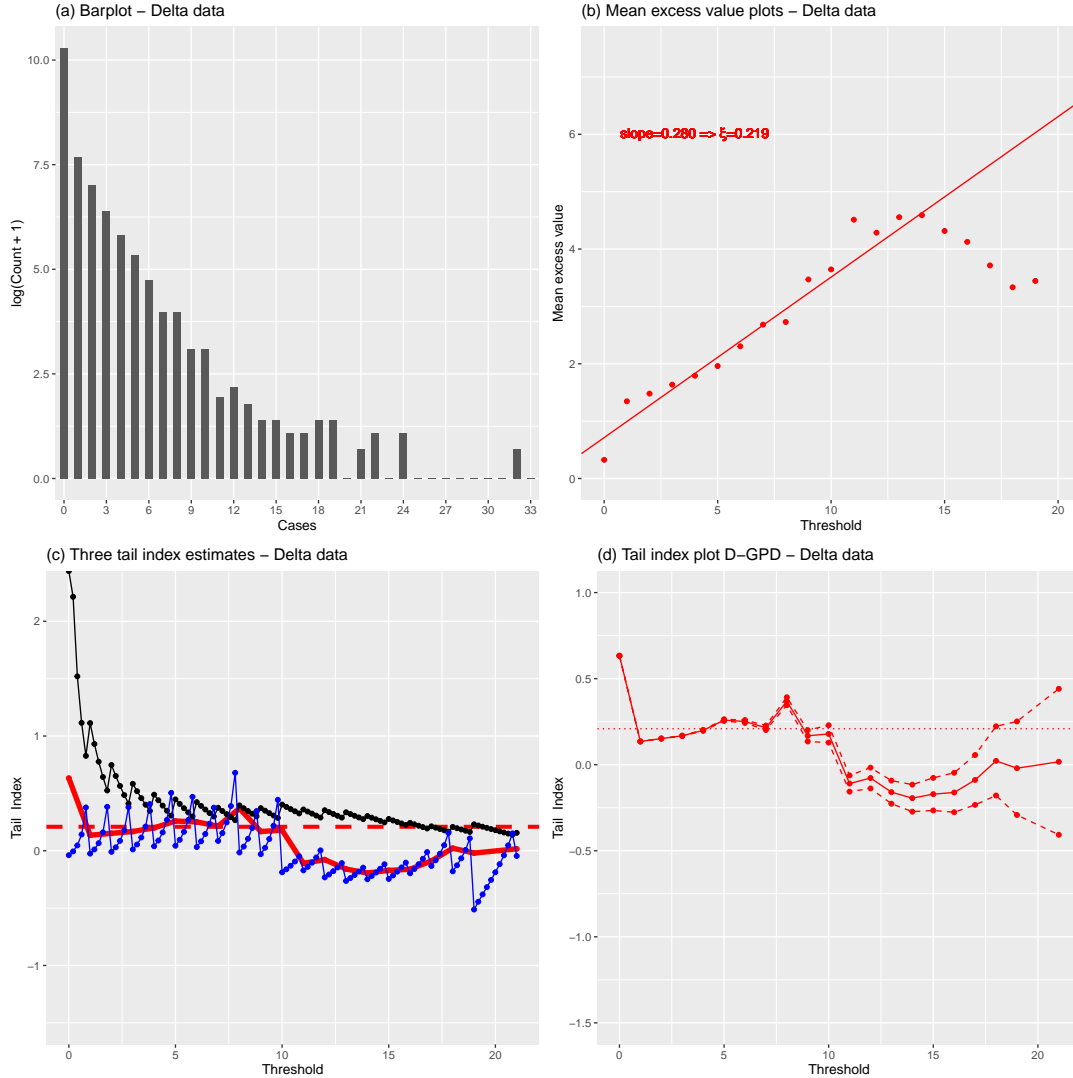


Figure 5: Secondary case data (Database S5) for SARS-CoV-2 from South Korea (July-August 2021). (a) Bar plot of the $\log(Z_i + 1)$ ($n = 33,903$). (b) Mean excess plots of secondary cases. (c) Hill (solid black), continuous GPD maximum likelihood (solid blue) and discrete GPD maximum likelihood (solid bold red) estimates of ξ . (d) Discrete GPD maximum likelihood estimates of ξ (solid red) and their associated 90% confidence intervals (dashed red). In panels (c) and (d), the averaged discrete GPD estimate $\hat{\xi} = 0.209$ over the stable region $u \in [1, 10]$ is indicated with the horizontal red line.

318 Californian data might be more conclusive. For the UK & Ireland dataset, the fact
319 that the sample is so small (30 clusters) in two countries with highly developed
320 healthcare and contact tracing systems is suspicious and may suggest reporting
321 issues.

322 Using the D-GPD model, one can gain further insight into large cluster sizes
323 by providing extrapolated estimates of extreme percentiles q_α potentially beyond
324 the sample maximum, through the estimate \hat{q}_α described in the Methods section.
325 Estimated 95th and 99th percentiles are given in Table 4. One may also match
326 the estimated percentiles with actual observations to get a sense of what would
327 constitute a conducive environment for the formation of large SARS-CoV-2 clusters.
328 For example, the estimated 95th percentile of 120 cases in Kerala is close to two
329 clusters of 113 cases (nursing home) and 132 cases (local transmission) already
330 observed in Kerala. Likewise, the estimate $\hat{q}_{0.95} = 272$ cases in Canada is fairly
331 close to a cluster of 324 cases in Canadian nursing homes. In Oregon, the estimated
332 99th percentile $\hat{q}_{0.99} = 124$ cases is in the vicinity of a cluster of 134 cases in a care
333 home setting. In Colorado, the estimate $\hat{q}_{0.99} = 140$ cases is close to a cluster of 134
334 cases in a nursing home. All of these clusters bar one (the local transmission cluster
335 in Kerala) correspond to indoor environments where social distancing is difficult to
336 practice.

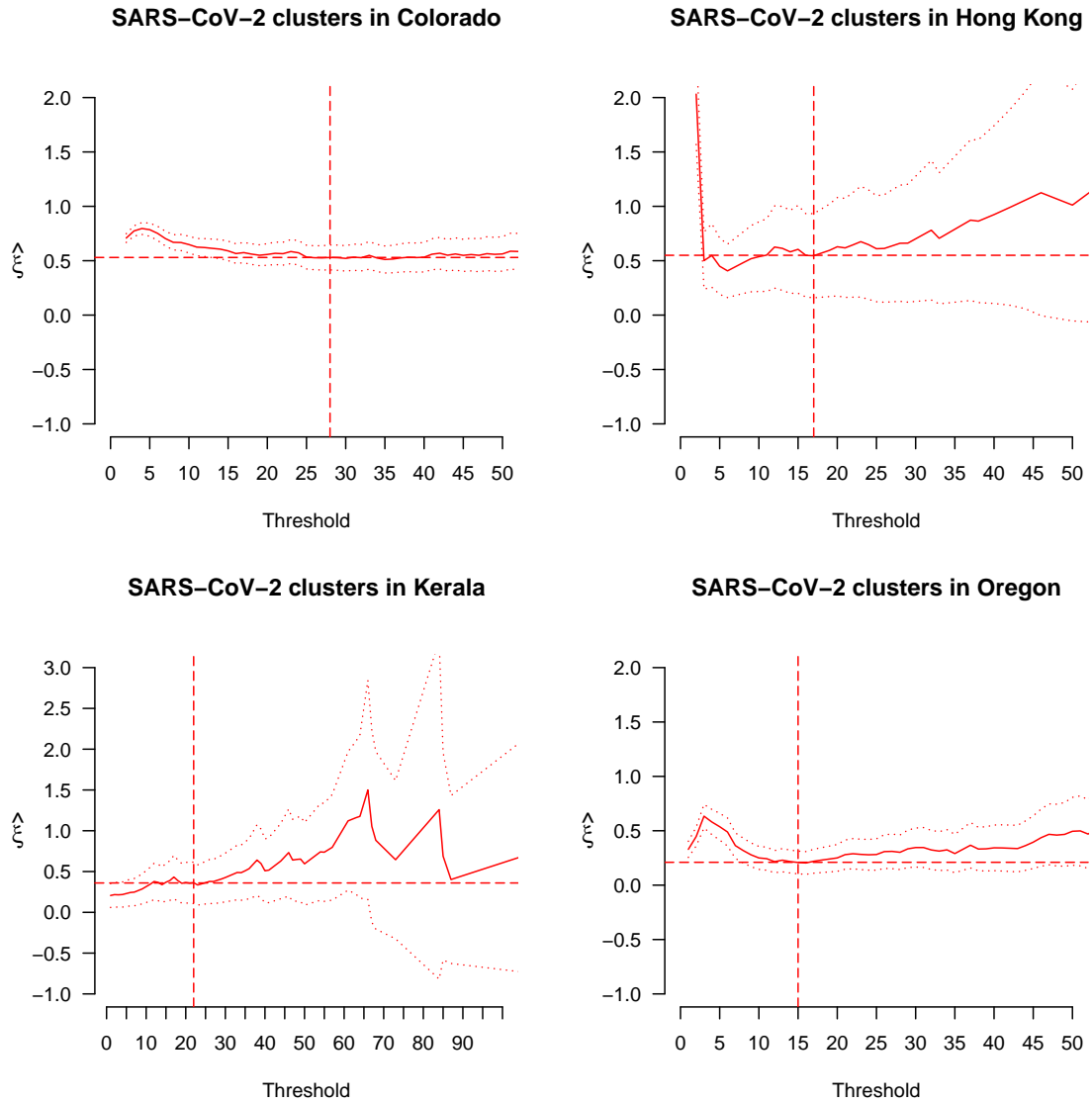


Figure 6: Analysis of cluster cases, for the four countries/states where the source is official data (Database S6). Plots of discrete GPD maximum likelihood estimates of ξ (solid lines), along with their 90% confidence intervals (dotted lines) and the final selected estimates (horizontal dashed lines) and thresholds (vertical dashed lines).

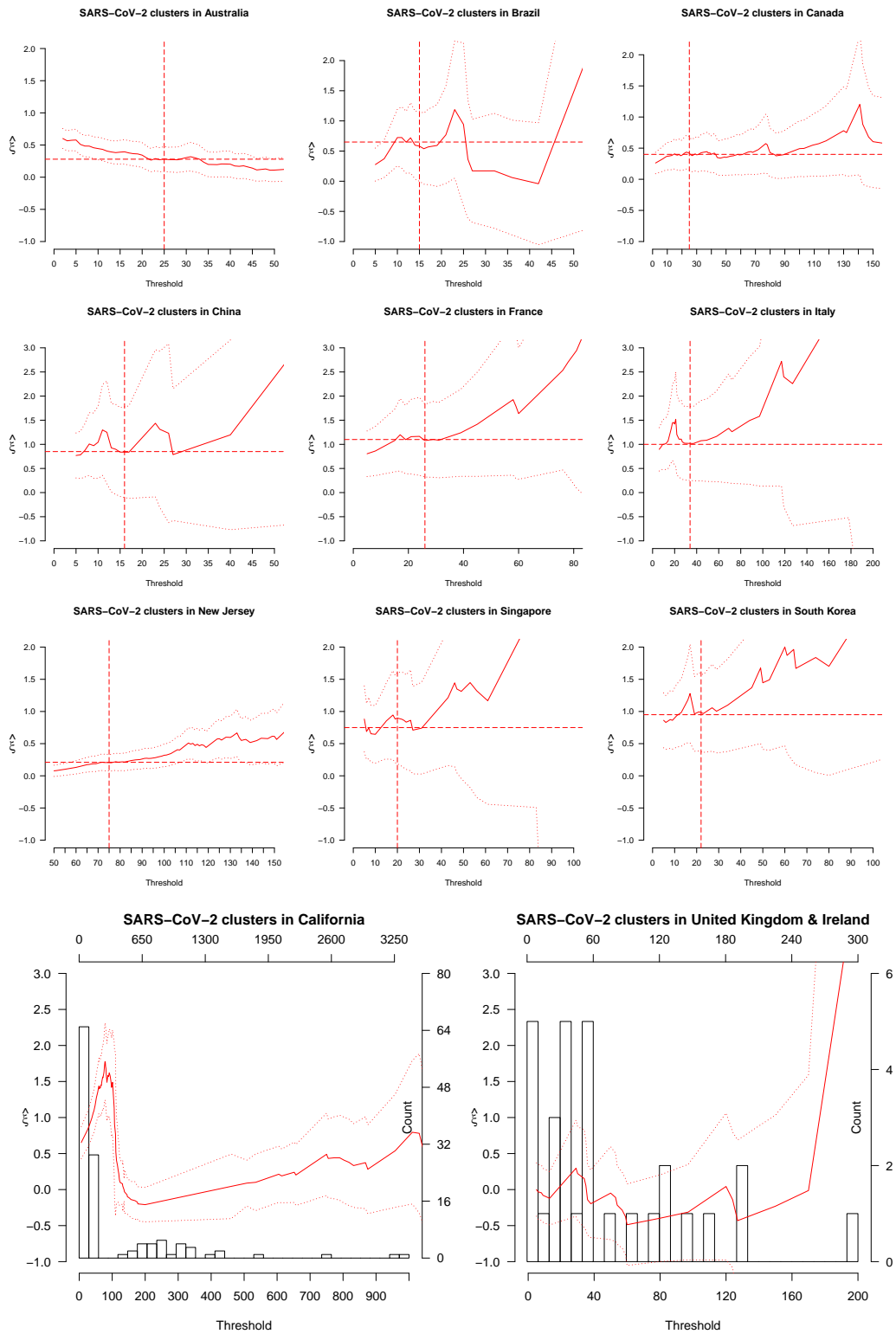


Figure 7: Analysis of cluster cases as in Fig. 6, with the results obtained from the data whose sources were the media (Database S7). The top 9 plots refer to those countries and states for which the extreme value analysis was conclusive. The bottom 2 plots refer to those for which the extreme value analysis was inconclusive.

Database S6

Location	n	$\hat{\xi}$ [90% CI]	u (n_u)	$\hat{q}_{0.95}$	$\hat{q}_{0.99}$	Max. Y_i (Setting)
Colorado, USA	4,769	0.53 [0.41, 0.64]	27 (474)	48	140	1,761 (Prison)
Hong Kong	54	0.55 [0.16, 0.93]	17 (34)	119	310	732 (Dancing)
Kerala, India	113	0.36 [0.11, 0.62]	22 (60)	120	255	580 (Unknown)
Oregon, USA	795	0.21 [0.10, 0.31]	15 (254)	64	124	639 (Prison)

Database S7

Location	n	$\hat{\xi}$ [90% CI]	u (n_u)	$\hat{q}_{0.95}$	$\hat{q}_{0.99}$	Max. Y_i (Setting)
Australia	355	0.28 [0.09, 0.48]	25 (145)	157	326	662 (Cruise ship)
Brazil	42	0.58 [0.00, 1.16]	15 (22)	82	220	191 (Hospital)
Canada	100	0.42 [0.15, 0.69]	25 (74)	272	624	1,500 (Meat processing plant)
China	34	0.84 [-0.12, 1.80]	16 (10)	99	401	368 (Market)
France	29	1.08 [0.32, 1.83]	26 (17)	443	2,530	2,500 (Religious gathering)
Italy	41	1.02 [0.25, 1.79]	34 (15)	378	2,013	7,000 (Stadium)
New Jersey, USA	183	0.20 [0.08, 0.33]	75 (157)	299	496	1,042 (Prison)
Singapore	45	0.90 [0.19, 1.61]	20 (21)	156	661	797 (Worker housing)
South Korea	45	0.98 [0.37, 1.59]	22 (24)	324	1,616	5,016 (Religious gathering)

Table 4: Final results for SARS-CoV-2 cluster sizes by country (first column), the corresponding sample size n (second column), D-GPD maximum likelihood ξ estimate and 90% confidence interval (third column), selected cluster size threshold u and associated number n_u of exceedance values $Y_i \geq u$ given $Y_i \geq u$ upon which the ξ estimate is calculated (fourth column), D-GPD maximum likelihood 95% and 99% percentile estimates of cluster size (fifth and sixth columns), and the sample maximum (last column). The top table corresponds to data from official sources (Database S6), and the bottom table to data from media sources (Database S7). The results reported in the latter table only concern the 9 countries and states for which the extreme value analysis was conclusive.

337 Discussion

338 In summary, we have investigated four datasets of secondary case numbers Z_i for
339 SARS-CoV-2 as a way to estimate and infer the extreme value index of the related
340 underlying offspring distribution. Motivated by the highly discrete nature of such
341 data, we used the Discrete GPD (D-GPD) maximum likelihood estimation method
342 which produces smoother and more stable plots of the associated D-GPD estimator
343 than the classical continuous GPD and Hill estimators. We first provided evidence
344 that the small SSE dataset (Dataset S2) compiled by [3] during the early phase of
345 the COVID-19 pandemic was fat-tailed, thus confirming their findings, although we
346 show in various ways that this dataset should not be pooled with their 15 SSEs
347 associated with SARS-CoV (Dataset S1), since they correspond to substantially
348 different distributions. On the other hand, as accurate extreme value inference
349 requires a large sample size in general, we also analysed an Indian secondary case
350 dataset of size 88,527 collected in 2020 (Database S3), which contains a very large
351 number of tied observations. The D-GPD estimate of the tail index is around 0.24,
352 which is in full agreement with the estimate of around 0.25 found by revisiting
353 the small SSE dataset of size 45 from [3]. The distribution of SARS-CoV-2 SSEs
354 therefore appears to have at least a finite third moment, whereas that of SARS-
355 CoV SSEs is found to have a much heavier upper tail with infinite variance and
356 therefore stronger superspreading effect. In an effort to account for the quality of
357 implemented control programmes as well as the nature of the variant under study,
358 we used two extra South Korean contact-tracing datasets. For the first dataset
359 (Database S4), collected in the first half of 2020 and used in [3], we cannot disprove
360 that the distribution of the number of secondary cases is light-tailed. By contrast,
361 for the second South Korean dataset (Database S5) collected during the summer of
362 2021, in which the majority of cases correspond to the Delta variant of SARS-CoV-
363 2 [5], we obtained a D-GPD estimate, $\hat{\xi} \approx 0.21$ clearly suggesting a heavier upper
364 tail for the Delta variant and therefore more pronounced superspreading potential
365 in South Korea relative to the first half of 2020.

366 We broaden our analysis by providing evidence that SARS-CoV-2 cluster sizes
367 are typically fat-tailed, based on 15 samples from 11 countries and 4 US states.
368 We infer the risk exposure and risk category of each country and state by making
369 use of D-GPD maximum likelihood estimates of both the extreme value index and
370 extreme percentiles, along with their associated confidence intervals. For the sake
371 of simplicity, we used a straightforward threshold selection rule, which is to spot a
372 stability region in the estimates (as a function of the threshold value) and choose an
373 estimate whose value is representative of those reached in this region. This practice,
374 colloquially known as “eyeballing”, is standard in applied extreme value analysis:
375 see for example the discussion in p.77 of Chapter 4 in [24]. It applies reasonably
376 well to the D-GPD sample paths, because they are overall much smoother and more
377 stable than the standard Hill and GPD maximum likelihood sample paths, which
378 are not designed to handle the discreteness of the data. The development of more
379 elaborate statistical techniques for the choice of threshold in discrete GPD maximum
380 likelihood estimation, such as methods based on asymptotic MSE minimisation or
381 the bootstrap in the spirit of the approaches outlined in Section 5.4 of [25] for Hill

382 estimation, is an open question which is beyond the scope of this paper.

383 A limitation of our study lies in the quality of the data, as it is not obvious
384 whether all SSEs or clusters over a given time period were available, or whether
385 cluster sizes were correctly recorded. To check robustness against missing data, we
386 have reproduced part of our analysis of cluster data by removing 10% of observa-
387 tions at random in each sample containing at least 100 data points, and replicating
388 this experiment 10,000 times. Robustness against poor recording was checked by
389 multiplying each observation Y_i by an independent normal variate W_i having mean
390 $\mu = 1$ and standard deviation $\sigma = 0.05$, and then reproducing our analysis of cluster
391 data on the $Y'_i = W_i Y_i$, this experiment being again replicated 10,000 times. There
392 is indeed some variation in the resulting estimates of ξ (Figs. 8 and 9), but this
393 does not affect our conclusion on the fat-tailed behaviour of the data, except in rare
394 situations when almost all the large values in the data go missing. This highlights
395 the importance of accurate data reporting as a prerequisite to such analyses.

396 It should be noted that, in classical epidemiological models, accurate estimation
397 of the basic reproduction number R_0 is of crucial importance as it informs the
398 extent of restrictions on social interactions and other control measures that should
399 be imposed to terminate the spread of an epidemic. The range of R_0 for SARS-CoV-
400 2 has been revised in [26] to 4.7-11.4, which is considerably higher than most early
401 estimates. This might explain why moderate restrictions that were implemented in
402 some nations, e.g. France, Italy, Spain, the UK, Australia and New Zealand, turned
403 out to be insufficient and replaced by nationwide or statewide lockdowns and/or
404 border closures. It should be clear that our results are, by construction, robust
405 to misspecified estimates of the expected number of secondary cases R_0 since they
406 solely rely on extreme values of numbers of secondary cases.

407 Our approach can be viewed as a proof of concept that transmission data from
408 a respiratory disease should not be pooled with data from a similar disease, since
409 similar R_0 numbers or parameters of average transmission do not, in general, induce
410 similar parameters of large community transmission. As such, preparing proactive
411 control measures actually requires a fine assessment of how unequal the distributions
412 of SSEs associated with different SARS-CoV-2 variants are. [27] conclude that the
413 reproductive number of the Delta variant is far higher than that of the historical
414 SARS-CoV-2 virus. Similarly, [28] estimate that the effective reproduction number
415 of the Omicron variant is more than 3 times that of the Delta variant in Denmark.
416 Our analysis of secondary case data did not, strictly speaking, allow one to conclude
417 statistically that SSEs linked to the Delta variant had a different extreme value index
418 from those linked to the original strains of SARS-CoV-2. However, in the contact-
419 tracing data recorded in South Korea, we did find a heavy tail in the offspring
420 distribution when the Delta variant made the majority of cases, as opposed to when
421 it did not. This tentative finding of a heavier tail in the data linked to the Delta
422 variant is coherent with the higher reproductive number of the Delta variant found
423 in [27]. The question of estimating parameters of large community transmission for
424 the Omicron variant remains open, as we could not find a dataset whose sample size
425 would enable us to draw statistically principled conclusions about the tail behaviour
426 of Omicron-related SSEs.

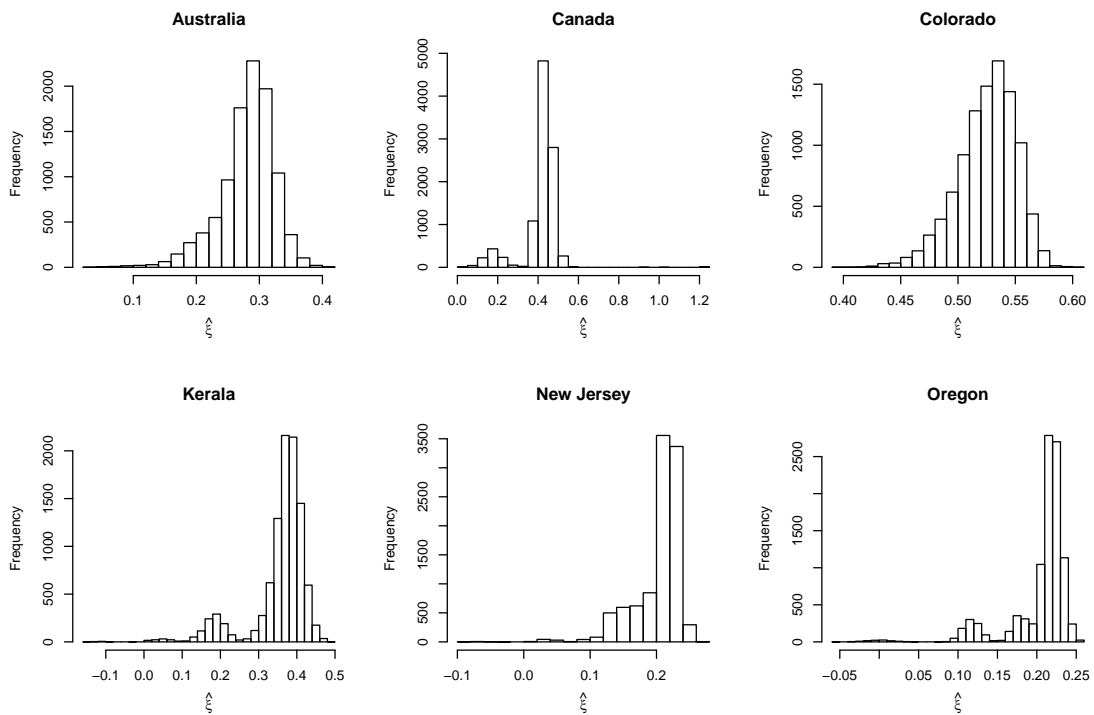


Figure 8: Robustness check (with respect to data omission) for the analysis of cluster cases (Databases S6 and S7). Histograms of the 10,000 estimates of ξ obtained by omitting at random 10% of the data. This was done only for the six samples containing at least 100 data points.

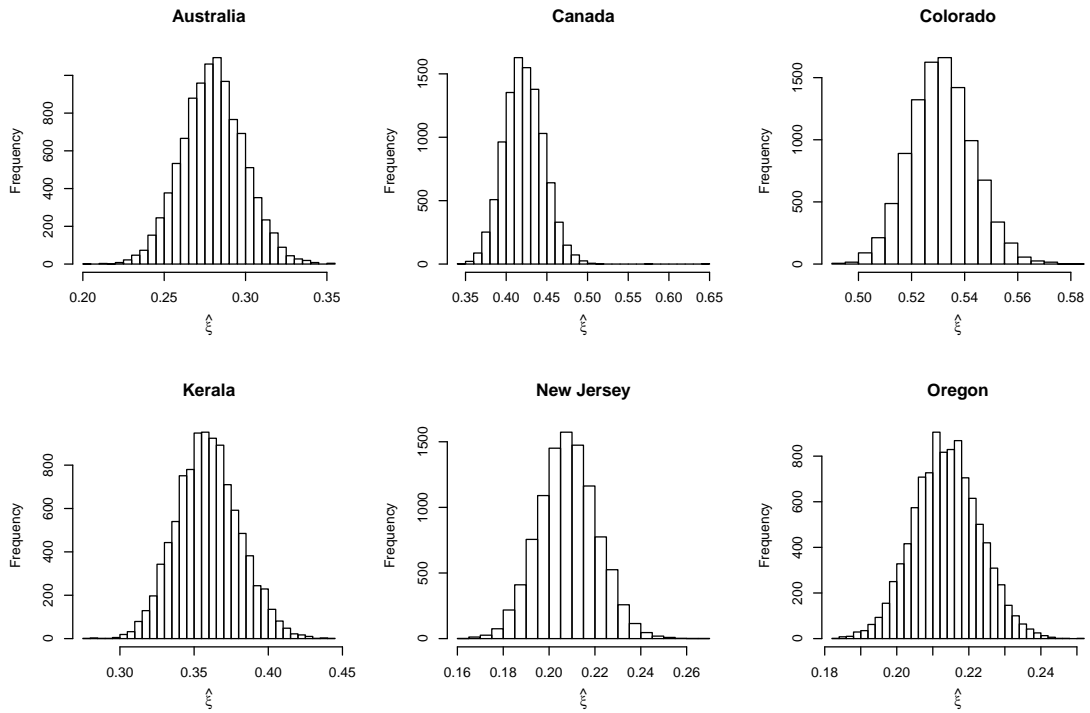


Figure 9: Robustness check (with respect to poor recording of the data) for the analysis of cluster cases (Databases S6 and S7). Histograms of the 10,000 estimates of ξ obtained by multiplying each data point by a random draw from the normal distribution with mean $\mu = 1$ and standard deviation $\sigma = 0.05$. This was done only for the six samples containing at least 100 data points.

427 **Ethics.** This article does not present research with ethical considerations.

428 **Data accessibility.** All datasets and the R code used for their statistical analysis
429 are available at <https://github.com/AntoineUC/SARS-CoV-2-codes>. Datasets
430 S1 and S2 can be found in the file `sse.R`. Database S3 can be found in the file
431 `traceDatSaved.Rdata`. Databases S4 and S5 can be found in the files `sse_korea_2020.txt`
432 and `sse_korea_2021_period_2.txt`, respectively. Databases S6 (apart from the
433 Colorado data) and S7 can be found in the file `clusters.Rdata`, while the Colorado
434 cluster size data can be found in `colorado.txt`.

435 **Authors' contributions.** A.U.C. undertook data curation and wrote the code for
436 the statistical analysis and visualisation of the results. All three authors participated
437 in the statistical analysis of the data and in drafting and revising the manuscript.

438 **Competing interests.** The authors declare no competing interests.

439 **Funding.** This research was supported by the French National Research Agency
440 (grant numbers ANR-19-CE40-0013, ANR-17-EURE-0010). G. Stupfler also ac-
441 knowledges support from an AXA Research Fund Award on 'Mitigating risk in the
442 wake of the COVID-19 pandemic'.

443 **Acknowledgements.** The authors acknowledge an anonymous Associate Editor
444 and two anonymous reviewers for their very helpful comments that led to a greatly
445 improved version of this paper.

446 References

- 447 [1] J. O. Lloyd-Smith, S. J. Schreiber, P. E. Kopp, W. M. Getz, Superspreading and
448 the effect of individual variation on disease emergence. *Nature* **438**, 355-359 (2005).
449 (doi:10.1038/nature04153)
- 450 [2] D. C. Adam, P. Wu, J. Y. Wong, E. H. Y. Lau, T. K. Tsang, S. Cauchemez, G. M. Leung,
451 B. J. Cowling, Clustering and superspreading potential of SARS-CoV-2 infections in Hong
452 Kong. *Nat. Med.* **26**, 1714-1719 (2020). (doi:10.1038/s41591-020-1092-0)
- 453 [3] F. Wong, J. J. Collins, Evidence that coronavirus superspreading is fat-tailed. *PNAS* **117**,
454 29416-29418 (2020). (doi:10.1073/pnas.2018490117)
- 455 [4] R. Laxminarayan, B. Wahl, S. R. Dudala, K. Gopal, C. Mohan B, S. Neelima, K. S. Jawa-
456 har Reddy, J. Radhakrishnan, J. A. Lewnard, Epidemiology and transmission dynamics of
457 COVID-19 in two Indian states. *Science* **370**, 691-697 (2020). (doi:10.1126/science.abd7672)
- 458 [5] S. Ryu, D. Kim, J.-S. Lim, S. T. Ali, B. J. Cowling, Serial interval and transmission dynamics
459 during SARS-CoV-2 Delta variant predominance, South Korea. *Emerg. Infect. Dis.* **28**, 407-
460 410 (2022). (doi:10.3201/eid2802.211774)
- 461 [6] State of Colorado, <https://covid19.colorado.gov/covid19-outbreak-data>, last updated
462 on 2 June 2021 (resolved outbreaks only). Accessed 27th September 2021.
- 463 [7] Government of Hong Kong, [https://www.chp.gov.hk/files/pdf/local_situation_](https://www.chp.gov.hk/files/pdf/local_situation_covid19_en.pdf)
464 [covid19_en.pdf](https://www.chp.gov.hk/files/pdf/local_situation_covid19_en.pdf), last updated on 6 September 2021. Accessed 6th September 2021.

- 465 [8] Government of Kerala, <https://covid19jagratha.kerala.nic.in/home/clusterList>. Ac-
466 cessed 21st July 2021.
- 467 [9] State of Oregon, [https://www.oregon.gov/oha/covid19/Documents/DataReports/
468 Weekly-Outbreak-COVID-19-Report-2021-08-25-FINAL.pdf?utm_medium=email&utm_
469 source=govdelivery](https://www.oregon.gov/oha/covid19/Documents/DataReports/Weekly-Outbreak-COVID-19-Report-2021-08-25-FINAL.pdf?utm_medium=email&utm_source=govdelivery). Accessed 25 August 2021.
- 470 [10] K. Swinkels, SARS-CoV-2 Superspreading Events Database, [https://kmswinkels.
471 medium.com/covid-19-superspreading-events-database-4c0a7aa2342b](https://kmswinkels.medium.com/covid-19-superspreading-events-database-4c0a7aa2342b). Accessed 21st
472 July 2021.
- 473 [11] T. Shimura, Discretization of distributions in the maximum domain of attraction. *Extremes*,
474 **15**, 299-317 (2012). (doi:10.1007/s10687-011-0137-7)
- 475 [12] F. Prieto, E. Gómez-Déniz, J. M. Sarabia, Modelling road accident blackspots data with the
476 discrete generalized Pareto distribution. *Accident Analysis & Prevention* **49**, 71:38 (2014).
477 (doi:10.1016/j.aap.2014.05.005)
- 478 [13] A. Hitz, R. Davis, G. Samorodnitsky, Discrete Extremes. arXiv [Preprint] (2017). [https:
479 //arxiv.org/abs/1707.05033](https://arxiv.org/abs/1707.05033) (doi:10.48550/arXiv.1707.05033)
- 480 [14] S. Ranjbar, E. Cantoni, V. Chavez-Demoulin, G. Marra, R. Radice, K. Jatton, Modelling the
481 extremes of seasonal viruses and hospital congestion: the example of flu in a Swiss hospital.
482 *J. Roy. Stat. Ser. C*, to appear (2022). (doi:10.1111/rssc.12559)
- 483 [15] L. de Haan, A. Ferreira, *Extreme Value Theory: An Introduction*, Springer-Verlag, New York
484 (2006). (doi:10.1007/0-387-34471-3)
- 485 [16] B. M. Hill, A simple general approach to inference about the tail of a distribution. *Ann.*
486 *Statist.* **3**, 1163-1174 (1975). (doi:10.1214/aos/1176343247)
- 487 [17] J. Beirlant, Y. Goegebeur, J. Segers, J. Teugels, *Statistics of Extremes: Theory and Applica-
488 tions*, John Wiley & Sons, Chichester (2004).
- 489 [18] S. Ghosh, S. Resnick, A discussion on mean excess plots. *Stoch. Proc. Appl.* **120**, 1492-1517
490 (2010). (doi:10.1016/j.spa.2010.04.002)
- 491 [19] S. Coles, *An Introduction to Statistical Modeling of Extreme Values*, Springer-Verlag, London
492 (2004). (doi:10.1007/978-1-4471-3675-0)
- 493 [20] P. Cirillo, N. N. Taleb, Tail risk of contagious diseases. *Nat. Phys.* **16**, 606-613 (2020).
494 (doi:10.1038/s41567-020-0921-x)
- 495 [21] C. Kremer, A. Torneri, S. Boesmans, H. Meuwissen, S. Verdonshot, K. Vanden Driessche,
496 C. L. Althaus, C. Faes, N. Hens, Quantifying superspreading for COVID-19 using Poisson
497 mixture distributions. *Sci. Rep.* **11**, 14107 (2021). (doi:10.1038/s41598-021-93578-x)
- 498 [22] N. Islam, Q. Bukhari, Y. Jameel, S. Shabnam, A. M. Erzurumluoglu, M. A. Siddique, J. M.
499 Massaro, R. B. D'Agostino, COVID-19 and climatic factors: A global analysis. *Environmental
500 Research* **193**, 110355 (2021). (doi:10.1016/j.envres.2020.110355)
- 501 [23] H. Hwang, J.-S. Lim, S.-A. Song, C. Achangwa, W. Sim, G. Kim, S. Ryu, Transmission
502 dynamics of the Delta variant of SARS-CoV-2 infections in South Korea. *J. Infect. Dis.* **225**,
503 793-799 (2022). (doi:10.1093/infdis/jiab586)
- 504 [24] M. Jacob, C. Neves, D. Vukadinović Greetham, *Extreme Value Statistics*, in: Forecasting and
505 Assessing Risk of Individual Electricity Peaks. Mathematics of Planet Earth, Springer, Cham
506 (2020). (doi:10.1007/978-3-030-28669-9)
- 507 [25] M. I. Gomes, A. Guillou, Extreme value theory and statistics of univariate extremes: A review.
508 *Int. Stat. Rev.* **83**, 263-292 (2015). (doi:10.1111/insr.12058)
- 509 [26] M. Kočańczyk, F. Grabowski, T. Lipniacki, Super-spreading events initiated the exponential
510 growth phase of COVID-19 with \mathcal{R}_0 higher than initially estimated. *Royal Society Open
511 Science* **7**, 200786 (2020). (doi:10.1098/rsos.200786)

- 512 [27] Y. Liu, J. Rocklöv, The reproductive number of the Delta variant of SARS-CoV-2 is far
513 higher compared to the ancestral SARS-CoV-2 virus. *J. Travel Med.* **28**, taab124 (2021).
514 (doi:10.1093/jtm/taab124)
- 515 [28] K. Ito, C. Piantham, H. Nishiura, Relative instantaneous reproduction number of Omicron
516 SARS-CoV-2 variant with respect to the Delta variant in Denmark. *J. Med. Virol.* **94**, 2265-
517 2268 (2022). (doi:10.1002/jmv.27560)