

Progress in
Physical Geography

**Fitting Limit Lines (Envelope Curves) to Spreads of
Geoenvironmental Data**

Journal:	<i>Progress in Physical Geography</i>
Manuscript ID	PPG-21-006.R2
Manuscript Type:	Main Article
Keywords:	Limit lines, Trimming method, Quantile regression, Non-parametric maximum likelihood methods, Envelope curves
Abstract:	<p>Geoscientists frequently are interested in defining the overall trend in x-y data clouds using techniques such as least squares regression. Yet often the sample data exhibits considerable spread of y-values for given x-values, which is itself of interest. In some cases the data may exhibit a distinct visual upper (or lower) 'limit' to a broad spread of y-values for a given x-value, defined by a marked reduction in concentration of y-values. As a function of x-value, the locus of this "limit" defines a "limit line", with no (or few) points lying above (or below) it. Despite numerous examples of such situations in geoscience, there has been little consideration within the general geoenvironmental literature of methods used to define limit lines (sometimes termed 'envelope curves' when they enclose all data of interest). In this work, methods to fit limit lines are reviewed. Many commonly applied methods are ad-hoc and statistically not well founded, often because the data sample available is small and noisy. Other methods are considered which correspond to specific statistical models offering more objective and reproducible estimation. The strengths and weaknesses of methods are considered by application to real geoscience data sets. Wider adoption of statistical models would enhance confidence in the utility of fitted limits and promote statistical developments in limit fitting methodologies which are likely to be transformative in the interpretation of limits. Supplements, a spreadsheet and references to software are provided for ready application by geoscientists.</p>

SCHOLARONE™
Manuscripts

Fitting Limit Lines (Envelope Curves) to Spreads of Geoenvironmental Data

Abstract

Geoscientists frequently are interested in defining the overall trend in x-y data clouds using techniques such as least squares regression. Yet often the sample data exhibits considerable spread of y-values for given x-values, which is itself of interest. In some cases the data may exhibit a distinct visual upper (or lower) 'limit' to a broad spread of y-values for a given x-value, defined by a marked reduction in concentration of y-values. As a function of x-value, the locus of this "limit" defines a "limit line", with no (or few) points lying above (or below) it. Despite numerous examples of such situations in geoscience, there has been little consideration within the general geoenvironmental literature of methods used to define limit lines (sometimes termed 'envelope curves' when they enclose all data of interest). In this work, methods to fit limit lines are reviewed. Many commonly applied methods are *ad-hoc* and statistically not well founded, often because the data sample available is small and noisy. Other methods are considered which correspond to specific statistical models offering more objective and reproducible estimation. The strengths and weaknesses of methods are considered by application to real geoscience data sets. Wider adoption of statistical models would enhance confidence in the utility of fitted limits and promote statistical developments in limit fitting methodologies which are likely to be transformative in the interpretation of limits. Supplements, a spreadsheet and references to software are provided for ready application by geoscientists.

Key words

Limit lines; Envelope curves; Trimming method; Quantile regression; Non-parametric maximum likelihood methods.

30 I Introduction

31 Ordinary least squares regression analysis commonly is used to define the statistical relationship
32 between one or more explanatory variables (X) and a response variable (Y). Where a relationship
33 exists, the trend can be linear or non-linear. Due to inherent instability in environmental systems,
34 the influence of additional unidentified explanatory variables, and the uncertainty in the
35 measurement procedures used to define x - y data pairs, usually there is considerable scatter in a
36 data plot of y -values on x -values. For ordinary least squares the uncertainty or randomness is
37 assumed to lie within the measurements of the dependent variable Y and not within the
38 independent variable X . Where uncertainty occurs in both X and Y other methods such as errors-
39 in-variables regression, total least squares regression and the reduced major axis method apply.
40 Herein, we restrict our attention largely to applications using or motivated by the ordinary least
41 squares method. The paper is written for non-specialists in statistical line fitting so supplements,
42 a spreadsheet and references to software are provided for some methods. However, users are
43 strongly recommended to seek the advice of professional statisticians in fitting any limit lines.

44 Often interest lies not with identifying the central trend to the x - y data, but with whether the x - y data
45 tend to indicate that maximum values of Y occur for given values of $X = x$. In similar vein, a
46 minimum limit may occur in some data sets. Below, mainly we explore the issue of defining the
47 trends in maxima, although the same procedures apply to defining minima. In the case where
48 maxima are expected or suspected to occur, identifying the trend line of maximal values of Y for
49 any given series of values of X become a focus of enquiry. Given sufficient maximum values of Y ,
50 a clear limit may be visually evident, with smaller values of Y defining scatter below the limit line.
51 More often, a limited sample size of x - y pairs means that there is no clearly defined limit although
52 one may be suspected to exist from the data scatter, or a limit can reasonably be assumed or is
53 known from theory. Limit lines also are referred to as envelope curves.

54 ***Overarching Objective of the Data Analysis***

55
56 Herein we review various methods that have been used to fit limit lines. Although sometimes theory
57 has informed the fitting of limit lines in the literature, oftentimes such consideration is lacking. The
58 researcher should consider what are the known or expected key characteristics of the expected
59 limit lines in terms of the likely effect on the decisions that might arise from the analysis. Thus, it
60 is beneficial if the form of the likely limit line can be specified or parameterised from theory. Where
61 theory is lacking, logical reasoning can be applied, informed by previous considerations of
62 empirical x - y data pairs similar to the target set of observations. These two approaches may
63 involve writing down the options for the form of the equations relating X and Y : e.g. $Y = f(X)$ and
64 considering the implications of fitting functions of different form. Rather than just utilizing the
65 existing data set, the simple procedure outlined above can assist in deciding where additional x - y
66 data points should be collected to improve understanding of the form of the limit line function and

1
2
3 67 the quality of the final fit. Knowledge of some or all these issues can make it easier to specify how
4 68 to estimate limit lines.

5
6 69 **Figure 1 near here**

7
8
9 70

10 71 **2 Approaches to limit line estimation – a statistician view**

11
12 72 This section seeks to provide an intuitive but rational framework within which the fitting of limit lines
13 73 can be discussed, motivated by elementary statistical thinking. Thereafter in Section 3, the relative
14 74 merits of different approaches to estimation of limit lines, known to be used by practitioners and
15 75 reported in the literature, are considered with respect to this framework.

16
17 76

18
19 77 It is assumed that the researcher has a data set or sample of pairs of points (x,y) , which *a priori* is
20 78 believed to be characterised by one or more defined limit lines. It is assumed that the existence
21 79 and characteristics of the limit lines are informed at least to some extent by the data. Typically, it
22 80 is assumed that given any value x of X , the corresponding values of Y are independently distributed.
23 81 Within our schema, methods for estimating limit lines can be considered to fall into four categories:
24 82 inspection, theory, joint statistical models and conditional statistical models, discussed in turn
25 83 below.

26
27 84

28 85 2.1 Inspection

29
30
31 86 Where the scatter of (x,y) data tend to define a boundary, the most frequently used approach is to
32 87 draw a line by eye: i) just outside of the data cloud, or ii) through selected data points along the
33 88 margin of the data cloud (e.g. a convex hull might be adopted). The nature of the line, for example,
34 89 linear or non-linear might be constrained by any known or expected theoretical or previous
35 90 empirical behaviour of the phenomenon.

36 37 91 2.2 Theoretical limit

38
39 92 In some situations, a theoretical function defining an expected limit line can be considered along
40 93 with the data plot and the relationship between this function and the empirical data can be
41 94 considered. Such an approach is related to defining tolerance limits or a specification, which can
42 95 be completely independent of the distribution of the plotted sample statistic.

43
44 96

45
46 97

47 98 2.3 Joint statistical models

48
49 99 Joint statistical models, like their conditional counterparts discussed in Section 2.4, are attractive
50 100 since they introduce a degree of objectivity into the estimation of limit lines (certainly in contrast to
51 **inspection**). The challenge is to specify the statistical model for the limit line in a manner such that

101 (a) the model can be estimated reasonably using a sample of data, and (b) observations for which
 102 modelling assumptions appear invalid can be identified using appropriate diagnostics, and the
 103 model rejected in favour of better-fitting alternatives.

104 Joint statistical modelling treats both X and Y variables as random (with upper-case letters used to
 105 indicate this) and seeks to estimate their joint distribution $f_{\{X,Y\}}(x,y)$. Limit lines might then be
 106 defined in terms of a contour in x - y with given statistical properties. For example, points on the
 107 contour might correspond to some fixed (low) probability density $f_{\{X,Y\}}(x,y) = p$; or the closed
 108 contour may define a region of x - y space with desired probability p (typically near unity). A simple
 109 example might be an ellipse of minimum enclosed area which encloses all the observations. See
 110 Ross et al (2019) for a discussion of contour construction in the context of environmental
 111 engineering. The portion of the contour corresponding to large y -values might be used as the limit
 112 line.

113 More generally appropriate models might be used to describe the marginal characteristics of
 114 variable X independently of the variable Y . Then, after marginal transformation to standard scale,
 115 a dependence or copula model (see Joe, 2014) could be used to describe the joint structure of the
 116 data on standard uniform margins.

117 The joint statistical model therefore can be rather complex. In contrast, conditional statistical
 118 models (discussed next) characterise the distribution of $Y|x$ for different fixed values x . Note the
 119 close relationship between joint and conditional distributions: for continuous random variables X
 120 and Y , e.g., we can write $f_{\{X,Y\}}(x,y) = f_{\{Y|X\}}(y|x)f_X(x)$, relating joint and conditional densities.

121

122 2.4 Conditional statistical models for $Y|x$

123

124 The data can be used to estimate a statistical model for Y given $X = x$. These models assume that
 125 the response is random or uncertain, whereas the value x of the explanatory variable is known and
 126 free of uncertainty. Note that more sophisticated approaches (e.g. hierarchical Bayesian inference)
 127 exist which build considerably on the basic conditional model structure considered here. There are
 128 many types of conditional model, as outlined in more detail below.

129

130 2.4.1 Linear regression

131 An initial assumption might be a **simple linear regression** relationship

132

$$Y = a + bx + \sigma\epsilon$$

133 between Y and x might apply. Here the intercept and slope parameters are a and b , σ is the
 134 measurement standard deviation and ϵ is a random variable with standard Gaussian
 135 distribution. Extensions to linear regression models, allowing for uncertain explanatory variables
 136 X also, known as errors-in-variables models, include total least squares. In cases where the

1
2
3 137 overall data spread in Y is not excessive relative to that in x , regression analysis can be used to
4 138 define the trend and confidence limits for $Y|X = x$ (henceforth $Y|x$ where possible for brevity) for
5 139 any value of x . A selected confidence limit can assist in positioning an appropriate limit line. This
6 140 approach applies in cases where the chosen confidence limit encloses all or most of the data points.
7
8 141 The linear regression can be refined in many ways to make it more suitable as a representation of
9 142 a limit line. These refined regression models are referred to further in Section 3.
10
11 143

14 144 2.4.2 Parametric model

15 145 Generalising linear regression, it might be assumed that the probability distribution of $Y|x$ is no
16 146 longer a Normal distribution, but rather some other distribution the parameters of which have
17 147 known functional forms in x . The objective of the data analysis is then to estimate these parameters
18 148 using techniques such as maximum likelihood estimation; Pawitan (2001) provides an excellent
19 149 introduction. The limit line for given x might then correspond to an extreme quantile of the
20 150 distribution $Y|x$ estimated under the parametric model. Coles (2001, Chapter 4) provides
21 151 illustrations using extreme value analysis.
22
23 152

27 153 2.4.3 Non-parametric model

29 154 Extending Section 2.4.3, there is no need to assume a parametric form for the parameters of the
30 155 distribution of $Y|x$, whilst seeking to estimate an extreme quantile of $Y|x$. Instead, we might assume
31 156 e.g. that the variation of these parameters with x can be described in terms of a linear combination
32 157 of basis functions (such as splines) defined on the domain of x . The model fitting would then
33 158 amount to estimating basis coefficients, and hence the specific form of parameter variation with x .
34 159 A popular approach in this situation is **quantile regression**, which estimates the quantile $Q(x)$ of
35 160 Y for given value of x with a specific non-exceedance probability τ . Koenker (2005) and Hao and
36 161 Naiman (2007) provide excellent introductions to the theory and applications of quantile regression.
37 162 A limit line might then correspond to $Q(x)$ as a function of x for an extreme non-exceedance
38 163 probability e.g. $\tau=0.95$.
39
40 164

46 165 2.4.4 Mixture model

47 166 Another approach which can be considered non-parametric is a **mixture model** for $Y|x$ (Maller *et al.*
48 167 *al.*, 1983; Kaiser *et al.*, 1994 in the geoenvironmental literature). Here, it is assumed that $Y|x$ is
49 168 drawn from one of a number of different linear regression models. The modelling task is to estimate
50 169 the parameters of all the regression models, and the probability that a given (x,y) pair in the data
51 170 was drawn from each of the linear regression models. An expectation–maximization (EM)
52 171 algorithm can be used to achieve maximum likelihood estimation. McLachlan *et al.* (2019) provide
53 172 a useful review of finite mixture modelling.
54
55 173

59 174 2.4.5 Conditional models for $Y|Y > u(x), x$

1
2
3 175 Because the focus of interest is in the largest values of Y for given x , it might be reasonable to
4 176 focus attention on a sub-set of the data for which $Y|x$ exceeds some threshold $u(x)$ (which might
5 177 itself be defined using quantile regression). In this case a local model can be fit to the sub-sample,
6 178 using any of the techniques mentioned in 3.1-3.4. One choice of parametric model with strong
7 179 asymptotic motivation might be an **extreme value model**, under which $Y|Y>u(x)$, x might follow a
8 180 generalised Pareto distribution with unknown shape and scale parameters to be estimated. A
9 181 shape parameter estimated to be negative would indicate the existence of a finite upper limit for
10 182 $Y|x$ which might be taken as the limit line. A positive shape parameter estimate would indicate that
11 183 no upper limit to the distribution of $Y|x$ exists; in this case, the limit line for $Y|x$ might be defined as
12 184 an extreme quantile of the distribution $Y|x$ estimated using the fitted parametric
13 185 model. Sophisticated applications of extreme value analysis are prevalent in some environmental
14 186 sciences, including hydrology; Coles (2001) provides an introduction.
15
16
17
18
19
20
21
22

23 188 **3 Approaches to limit line estimation – a practitioner view**

24 189 This section lists some of the methods used by practitioners, and reported in the literature, for
25 190 estimation of limit lines. With reference to Section 2, this section also provides an outline of the
26 191 strengths and weakness of the various approaches. Methods are listed in approximate order of
27 192 increasing complexity.

28 193 **Inspection** (see Section 2.1) fits a line that often is referred to as an envelope curve and can 'over-
29 194 predict' the limit line if the line is drawn such that all data points lie below it. The 'true' limit line
30 195 could lie closer to the data than it is actually drawn. In this case no data points actually occur at
31 196 the limit – which is counterintuitive. The method has the advantage that eye-defined complex limits
32 197 can be drawn which might be difficult to define mathematically, or which might lack theoretical
33 198 justification. This latter advantage also can be considered a disadvantage, as subjectivity is
34 199 involved in positioning the line. If the purpose of fitting the line is merely to draw attention to the
35 200 possible presence of a limit then inspection is useful but it lacks objectivity. Examples of this kind
36 201 abound in the literature: for example, Innes (1983) fitted curves through the outermost data points
37 202 to define empirical lichen growth curves.

38 203 **Theoretical limit** (see Section 2.2) is a powerful means to define limiting lines. Theoretical curves
39 204 can be added to a graph without consideration of the empirical data, in which case the method
40 205 cannot be considered a fitting procedure. However, oftentimes theoretical curve fitting makes use
41 206 of the empirical data and so can be regarded as a fitting procedure. The relationship between the
42 207 trend of the theoretical curve, the position of individual data points, the configuration of clouds of
43 208 related points and the relative plotting positions of clouds can result in reflection as to the accuracy
44 209 of the individual data point values, the relationship between clouds, or consideration as to whether
45 210 the theory needs revision. Fitting of a theoretical curve, independently of any consideration of the
46 211 empirical data, can be epitomized by the classic concept of bedload transport efficiency (Bagnold,

1
2
3 212 1966) whereby Bagnold (1980; see also Carling 1985) compared empirical data with an efficiency
4 213 maximum function that effectively constitutes a limit line. In contrast, Kaiser *et al.*, (1994) used
5 214 ecological theory of limiting factors, informed by empirical data, to develop several statistical
6 215 approaches to fit limit lines to limnic biological process data. A worked example is provided in
7 216 section 5 and within Supplement 1.
8
9
10

11 217

12
13 218 **Environmental contours** (see e.g. Ross et al 2020, and Section 2.3) are popular in coastal and
14 219 offshore engineering. For two random variables X and Y , a closed contour is sought which encloses
15 220 a subset of the domain of the random variables with a given probability p just below unity. Regions
16 221 outside the subset are considered rare or extreme. The contour line itself can also be considered
17 222 a limit line. The location of the contour typically requires that a joint model for variables X and Y is
18 223 established. Extreme value analysis (see Section 2.4.5) is often an important ingredient in the
19 224 estimation of the joint model.
20
21
22
23
24

25 225 In **selective regression**, the limit line might be defined using a prior linear regression $y = a + b x$
26 226 through the whole sample (see Section 2.4.1). The limit line would also be linear in x , with the
27 227 same slope b as the linear regression line, and an increased value of intercept a^* , such that $y =$
28 228 $a^* + b x$ forms the appropriate limit line. A linear limit line located in this way is referred to as
29 229 selective regression, because it can be used to exploit knowledge of just some of the (x,y)
30 230 observations in the sample for analysis. We might consider fitting a linear regression (with fixed
31 231 slope b from the whole-sample regression) to a selected sub-sample of large values of $Y|x$ for
32 232 different x , as a more systematic procedure to estimate a^* . Because confidence limits for linear
33 233 functions are non-linear, the analyst might also exploit knowledge of confidence limits from a
34 234 whole-sample regression to select an appropriate value of a^* in selective regression, such that
35 235 the limit line is equivalent roughly to the selected confidence limit. Such an approach is similar to
36 236 the concept of applying 'control limits' also known as 'natural process limits' used in system
37 237 monitoring where, if there are sufficient normally distributed values of Y for a given value of x , a
38 238 limit is placed at a distance of ± 3 standard deviations (SD) from empirical estimate for the mean of
39 239 $Y|x$. For normally distributed values of Y for given x , 99.73% of all the plot points on the chart will
40 240 fall within the ± 3 SD limit. Thus only 0.27% of data points should lie above the limit line.
41
42
43
44
45
46
47
48
49

50 241 In selective linear regression, a whole-sample simple linear regression can be used to inform the
51 242 location of the limit line. The draw-back is that it can be difficult to select which data points should
52 243 be considered relevant for the specification of a new intercept a^* , especially difficult where the
53 244 data spread poorly defines a limit and where outliers are frequent. Selection of the points used to
54 245 define the limit is largely subjective. In the example (Fig. 1A), fortunately there are no distinct
55 246 outliers and the regression lines were fitted through an eye-selected set of 'outer' points. In this
56 247 example, the procedure leaves no points above the limit lines, but where outliers exist the
57
58
59
60

248 procedure is clearly unsatisfactory unless the outliers are excluded objectively. Thus, assessing
249 the influence of outer points can assist in the decision making (see below). In selective linear
250 regression, the offset limit curve can be assumed to have the same form as the least squares
251 regression function fitted to all the **data; however, this may not always be the case**. Clearly, the
252 method cannot apply when the spread of the data visually indicates that the limit line does not have
253 the same trend as the least-squares model applied to the complete sample (e.g. Fig. 1A). In many
254 cases the data spread is not considered and the analyst simply fits a curve parallel to the least-
255 squares fit to all data (Fig. 1B). However, if there are sufficient normally distributed data for Y given
256 x , then placing a limit at ± 3 standard deviations (SD) from the central tendency of the trend is
257 rational and reproducible. As examples, Gaume *et al.*, (2009) and Tarolli *et al.* (2007) fit limit lines
258 to extreme flood data (flood envelope curves) using selective regression, whilst Castellarin (2007)
259 briefly reviews the history of this approach to the development of flood envelope curves, and
260 introduces a probabilistic method to consider the likelihood of floods exceeding the limit curves.
261 More robust statistical methods are preferable, including linear quantile regression (Fig. 1C). In
262 this example, visual inspection indicates that there is a significant number of potential outliers
263 above the 0.90 quantile, in contrast to the situation within Fig. 1 A & B. So, in the case of Fig. 1C,
264 identification of the appropriate quantile and identifying outliers needs addressing further. The
265 example in Fig. 1C is considered again below.

266
267 The **iterative selective regression** procedure of Maller *et al.* (1983) is an iterative least squares
268 procedure in which data points are down-weighted according to their distance from a trial line to
269 obtain a new line. This latter line forms the basis for the next iteration. This procedure is equivalent
270 to fitting the least squares line through an objectively-derived subset of the data. For consistency
271 with our notation, we refer to this approach as iterative selective regression, although Maller *et al.*
272 (1983) referred to it as a trimming method. Simulations of the estimates for the iterative selective
273 regression approach show that small biases occur, but the estimates of slope and intercept are
274 approximately normally distributed and are reproducible by other operators. The solution is not
275 uniquely determined, but the accepted fitted line usually is taken to be the solution that includes
276 the greatest number of data points. Carling (1987) used the Maller *et al.*, (1983) method to fit a
277 limit line to define a maximum lichen growth curve. Guidance notes on implementing the Maller *et al.*
278 *et al.* (1983) method and an Excel work sheet are archived on Github (Carling *et al.*, 2021).

279
280 **Parametric model fitting** (see Section 2.4.2) is widespread in environmental sciences. Once the
281 parameters of the model have been estimated by fitting to the complete sample, the limit line can
282 be specified and easily calculated e.g. in terms of a quantile of the conditional distribution $Y|X = x$.
283 Fundamental physical and statistical considerations often motivate the choice of parametric model.
284 For example, for count data a Poisson model might be appropriate (see e.g. Chavez-Demoulin
285 and Davison 2005). For measurements of contaminant levels in soils, a log-normal or gamma

1
2
3 286 distribution is often appropriate. The simple linear regression model of Section 2.4.1 is an example
4
5 287 of parametric model fitting using a Gaussian assumption. Polynomial models of the form e.g.
6
7 288 $Y|(X = x) = ax + bx^2 + cx^3 + \sigma\epsilon$, and response surface models of the form e.g. $Y|(X_1 = x_1, X_2 = x_2$
8
9 289 $) = ax_1 + bx_2 + cx_1^2 + dx_2^2 + ex_{12} + \sigma\epsilon$ (in terms of two covariates X_1 and X_2) are also examples of
10
11 290 parametric models suitable to define limit lines. Davison and Ramesh (2000), Hall and Tajvidi
12
13 291 (2000) and Ramesh and Davison (2002) developed local likelihood models for smoothing sample
14
15 292 extremes of single series. Response surface methodology (RSM) is a tool that was introduced in
16
17 293 the early 1950s by Box and Wilson (1951). RSM is a collection of mathematical and statistical
18
19 294 techniques that is useful for the approximation and optimization of multivariate stochastic models
20
21 295 of 3D surfaces. For example, Shirazi *et al.* (2020) applied RSM techniques to multivariate data to
22
23 296 fit optimal maximum response surfaces related to factors controlling soil erosion using an objective
24
25 297 function they termed the desirability function.

26
27 298
28
29 299 Eberhardt and Thomas (1991), considering environmental systems, recommend the **Box and**
30
31 300 **Lucas method** to obtain optimal parameter estimates of response surfaces; thus effectively
32
33 301 defining limit lines. Box and Lucas is a relatively robust approach but implementation needs a
34
35 302 higher level of statistical competency, although software is available to fit a selection of functions
36
37 303 (e.g. Originlab®). The original use was to define a complex curve through few data points which
38
39 304 are believed to be the optimal (or in our case maximal) values of Y for given x , to thus assist in
40
41 305 choosing further values of x to sample for Y . As new data are added the line is optimized again.
42
43 306 The procedure assumes that the trend of the final fitted line defines the outer limit of the region
44
45 307 within which data might be expected to occur, or which points are operationally acceptable. Thus,
46
47 308 the method is heavily dependent on some prior knowledge of the expected behaviour of maximal
48
49 309 values of Y as a function of x . Box and Lucas (1959) did not consider the case where there are
50
51 310 many sub-optimal values of Y , which is the focus of this paper. Consequently, there is an issue as
52
53 311 to the initial section of points for fitting in cases where many sub-optimal values of Y exist.

54
55 312 **Quantile regression** (see Section 2.4.3) is capable of modelling any specific quantile of the
56
57 313 conditional distribution $Y|X = x$ including the tails (corresponding to say the 95% quantile).
58
59 314 However, good performance requires sufficient data to characterise $Y|x$ reasonably as a function
60
315 of x . To estimate the 95% quantile we therefore need considerably more than $1/(1-0.95)=20$
316 observations of Y in the vicinity of each value of x of interest; for the 99% quantile, in excess of
317 100 observations are required locally for each x . Compared with linear regression, quantile
318 regression is computationally somewhat more demanding, and typically performed using software
319 such as R, PYTHON or MATLAB. Extensions of quantile regression to estimate non-crossing
320 quantiles simultaneously corresponding to different non-exceedance probabilities are
321 computationally more demanding still. Cade (2017) provides an outline of the method for
322 environmental sciences. A simple example is presented as Fig. 1C and a further example is
323 provided in Supplement 2.

1
2
3 324 **The mixture model** of Maller *et al.*, (1983; Kaiser *et al.*, 1994, and see Section 2.4.4) needs a
4
5 325 reasonably high level of statistical competency. In outline, values of γ are assumed to be drawn
6
7 326 from a mixture of Gaussian distributions. The mean and standard deviation of each mixture
8
9 327 component is linearly related to a 'fullness' random variable drawn from $[0,1]$. The mean of each
10
11 328 mixture component is also related to x by a linear regression. During inference, a set number of
12
13 329 'fullness' values is considered, and the parameters of the linear regression and the mixture
14
15 330 component from which each particular pair (x,y) are drawn are estimated. The final choice of limit
16
17 331 line to adopt given the inference is the choice of the investigator.

18
19 332
20 333 **Extreme value analysis** (see Section 2.4.5) is used widely in environmental science to define
21
22 334 return values for processes such as rainfall, temperature, storm, wildfire and earthquake severity,
23
24 335 extreme occurrences of which are hazardous. The T -year return value is defined by the equation P
25
26 336 $(Y_A > y) = 1/T$, where Y_A is the annual maximum of random variable Y . The distribution of Y_A is
27
28 337 estimated based on a sample of data using extreme value analysis (see *e.g.* Coles 2001). The
29
30 338 return value can be also be defined conditional on a co-variate X , as $P(Y_A > y|X = x) = 1/T$. In
31
32 339 this case, a different return value is estimated for each value of x of X (see *e.g.* Davison and Smith
33
34 340 1990). Further details and a software reference are found in Supplement 3.

35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

342 **4. Practical issues**

343 A number of practical issues arise in attempting to estimate limit lines from a sample of data. In
344 this section, we provide an overview of some of the issues that are likely to be of concern to the
345 practitioner. These include identification of outliers, breakpoints and mixed samples, and the
346 quantification of uncertainty of inference.

347 **4.1 Identifying outliers**

348 In regression modelling (Section 2.4), observations with large residuals (outliers) or high leverage
349 are problematic, since they may violate the assumptions underlying the model and cast doubt on
350 the outcome of a regression. Outlier detection and regression diagnostics naturally have a large
351 statistical literature; the works of Wetherill *et al.* (1982) and Cook and Weisberg (1982) provide
352 introductions. Traditionally when assessing a dataset before conducting linear regression, outliers
353 were identified by eye from inspection of the x - y scatterplots. Objectively identified outliers likely
354 lie above any proposed limit line so their identification is critical when fitting limit lines.

355 Unusually large values of γ and x can be identified by examination of extreme quantiles of marginal
356 statistics. Alternatively, if sufficient data for γ occur for a given x , or within some neighbourhood
357 of x , then outliers can be identified from examining histograms of $\gamma|x$ for each x of interest. Within
358 a linear regression context, model diagnostics such as the diagonal elements of the so-called hat
359 matrix, and Cook's distance can be used to identify observations with high leverage and influence

1
2
3 360 respectively. Large values of model fit residuals are indicative of outliers. It is also generally useful
4 361 to examine so-called studentized residuals. These diagnostics are explained in many statistical
5 362 treatises, and are often included in statistical software packages, so we do not elaborate further.

6
7
8 363 For joint modelling of bivariate data (Section 2.3), Mahalanobis distance and similar metrics can
9 364 be used to identify data points which are unusual with respect to that metric. In a regression
10 365 context, when the occurrence of outliers can be attributed to one or more additional data-
11 366 generating processes (over and above those responsible for the bulk of the sample), then more
12 367 sophisticated techniques including mixture modelling can be used to simultaneously estimate 'bulk'
13 368 and 'outliers' (e.g. Aitkin and Tunnicliffe Wilson 1980, Yu *et al.* 2015).

17 369 **4.2 Identifying breakpoints**

18
19
20 370 It is possible that an attempt to estimate a limit line with given characteristics (e.g. linearity) through
21 371 x-y data does not yield satisfactory results. If the limit line is estimated using a statistical procedure,
22 372 then lack of fit can be quantified. In such cases, more general models for the limit line should be
23 373 sought. The relative performance of different models for the limit line can then be compared, and
24 374 the best model adopted (e.g. Wetherill *et al.* 1986). Sometimes it may be appropriate to consider:
25 375 (1) whether the data might exhibit breakpoints or changepoints in the x-y relationship, or; (2)
26 376 whether a model admitting a non-linear relationship between variates is appropriate (e.g. Zanini *et*
27 377 *al.* 2020). Figure 1C illustrates this issue. Here, the slope of the limit line clearly changes at wind
28 378 speed around 10 ms⁻¹; it might therefore be appropriate to fit a piecewise linear limit line as
29 379 illustrated. However, physically we know that water waves are generated by the wind via frictional
30 380 drag forcing, which implies alternative approaches including a linear limit line for y on the square
31 381 of x, or a quadratic quantile regression limit line might be appropriate. However, the relationship
32 382 observed at a specific location is unlikely to follow the quadratic form exactly, due to various effects
33 383 including fetch-limitation, wind-field non-stationarity, bathymetric effects in shallower water etc. For
34 384 this reason, fitting a piecewise linear form for a limit line is a pragmatically sound way to proceed;
35 385 in practice, a larger number of piecewise segments might probably be used. In fact, exactly this
36 386 approach is frequently used in ocean engineering to specify an extreme value threshold, and
37 387 amounts to an approximate non-parametric quantile regression (see Section 2.4.3).

38
39
40
41
42
43
44
45
46
47 388 In general, identifying breakpoints or changepoints in a sample can be important in the
48 389 interpretation of a physical process (e.g. Ryan *et al.*, 2002). The modelling challenge is to identify
49 390 one or more breakpoints in x in the sample such that limit lines using data in each of the resulting
50 391 sub-sets can be estimated more parsimoniously than using the complete sample. Often, prior
51 392 empirical knowledge, or theory, can be used to locate the breakpoints in terms of x-values. Then
52 393 separate regression models (or other approaches from Section 2) might be adopted for each sub-
53 394 set to estimate limit lines. When the location of a breakpoint is uncertain, data points close to the
54 395 expected breakpoint first can be considered to fall into one group and then considered to be part

396 of the other group to influence the regression line trends, thus repositioning the expected
397 breakpoint.

398 Identification of breakpoints also can be achieved as part of the statistical inference. For example,
399 optimal partitioning of the x -domain into K intervals, on each of which piecewise constant or
400 piecewise linear regression models are estimated, can be performed (see Ryan *et al.* 2002, Yang
401 *et al.* 2016).

402 **4.3 Identifying mixed samples**

403 Sometimes, it is possible that the sample for analysis corresponds to observations of a mixture of
404 different data-generating processes. In this situation, we might expect that limit lines would be
405 more appropriately estimated for the individual processes from which the mixture is composed. It
406 might therefore be useful to perform prior partitioning of the sample into two or more groups using
407 data for both Y and X . This outcome can be accomplished using cluster analysis when there is no
408 prior knowledge of group membership or using one of a large variety of classification techniques
409 (including random forests and support vector machines) when some knowledge of group
410 membership is available. For the two-group case, discriminant analysis (Brereton, 2009; Dixon and
411 Brereton, 2009) is another popular choice. As mentioned in the context of outlier detection above,
412 more sophisticated statistical techniques to model the mixture explicitly can also be employed.

413 **4.4 Quantifying uncertainty**

414 Quantifying the uncertainty of estimates of limit lines is generally important if those estimates are
415 to be trusted. Some of the approaches described in Sections 2 and 3 do not involve an explicit
416 quantitative model for the relationship between Y and X ; it is difficult therefore to quantify the
417 uncertainty with which these limit lines are estimated. Other methods from Section 2 and 3 make
418 combined use of a data sample and a statistical model; for these methods and the limit lines they
419 produce, it is therefore possible to quantify uncertainty using well-established approaches.

420 Sources of model uncertainty can be considered aleatory (due to the inherent natural variation of
421 the process we are modelling, which will always be present) or epistemic (due to inadequate data,
422 measurement procedures, model specification etc., the effects of which we could in principle
423 eliminate with enough effort).

424 When a regression-type model for $Y|X = x$ is being estimated, there are broadly two approaches
425 to the quantification of uncertainty. The first approach adopts Bayesian inference. The key steps
426 are (a) specification of full probabilistic data-generating model, (b) specification of a joint prior
427 distribution for all the parameters in the model, (c) estimation of the joint posterior distribution of all
428 parameters by conditioning on a sample of data using Bayes theorem, (d) diagnosing model
429 performance, and estimation of posterior predictive credible intervals for structure variables of
430 interest, such as a limit line. Many statisticians view Bayesian inference as the preferred strategy
431 for model building and decision making, but it often suffers because of the difficulty of specifying

1
2
3 432 reasonable prior distributions for parameters, and the computational complexity of inference.
4 433 Bishop (2006) and Gelman *et al.* (2015) provide introductions.

6 434 The second approach to uncertainty quantification is based on assessing the variability of
7 435 inferences from models estimated using resamples of the original data sample. Different
8 436 resampling techniques, including cross-validation, bootstrapping and randomised permutation
9 437 testing provide relatively simple pragmatic approaches to estimate the performance of statistical
10 438 models, to estimate uncertainties of predictions, and perform significance tests. Resampling
11 439 approaches are widespread in the applied literature, especially when there is some ambiguity
12 440 about the appropriateness of the model being used. However some might claim that resampling
13 441 approaches lack the overall coherence and elegance provided by the Bayesian approach. There
14 442 is a huge literature on resampling methods; Good (2006) provides an introduction. The works of
15 443 Molinaro *et al.* (2005), Hesterberg (2015) and Lehr and Ohm (2017) provide useful practitioner
16 444 perspectives.

23 445 **4.5 Measurement error and heteroscedasticity**

24 446 In many data sets, measurements of both Y and X are made with error. That is, we cannot measure
25 447 either Y or X precisely. Uncertainty in Y can be accommodated relatively easily in the distributional
26 448 assumption made for $Y|X = x$. However, uncertainties in X are more problematic to handle
27 449 appropriately in simple statistical models. The presence of measurement errors causes increased
28 450 bias and uncertainty in fitted statistical models, leading to erroneous inferences about limit lines.
29 451 Using Bayesian inference, we can routinely specify a measurement error model for both Y and X .
30 452 Alternatively, we can extend conventional regression models to so-called errors-in-variables
31 453 models.

32 454 In a simple linear regression model, we make the assumption that the variance of $Y|X = x$ does
33 455 not change with the value of X . However, in many applications, this is not the case, and the data
34 456 are said to exhibit heteroscedasticity. This feature can again be accommodated by extending the
35 457 regression model.

36 458 **4.6 Model selection**

37 459 Model selection is a procedure to select one among many candidate models. Typically we select
38 460 a model with the best performance for the task at hand. However, there may be many competing
39 461 issues relevant for good model selection other than quantitative performance, such as model
40 462 complexity and interpretability. In many practical situations, a model which is straightforward to
41 463 estimate, interpretable and gives reasonable performance, is preferable over a considerably more
42 464 complex model which is less interpretable and gives only slightly improved performance.

1
2
3 465 There are essentially two approaches to model selection. In general, probably the wisest approach
4 466 is based on the assessment of **predictive performance** of the model, preferring the candidate
5 467 model with best predictive performance. Predictive performance is assessed by quantifying out of
6 468 sample error; that is, how well a model performs on data that were not used to fit the model in the
7 469 first place. There are many approaches to quantifying predictive performance, including (1)
8 470 partitioning the original data into two groups, using one group to fit a model, and the other group
9 471 as an unseen test set to estimate predictive performance, and (2) cross-validation, in which the
10 472 original sample is partitioned into a number of subsets which are withheld one at a time, serving
11 473 as test sets for models estimated using all the remaining sets; an estimate of predictive
12 474 performance is then accumulated over all the test sets. The second approach to model selection
13 475 attempts to quantify model performance using **fitting performance** of the model. However,
14 476 because fitting performance is typically an over-optimistic assessment of predictive performance,
15 477 the fitting performance score is usually penalised by a measure of model complexity; more complex
16 478 models receive higher penalties. A number of related performance measures, including the Akaike
17 479 Information Criterion (AIC), Bayesian Information Criterion (BIC), and Minimum Description Length
18 480 (MDL) are available. Pawitan (2001, Sections 13.5-13.6), Davison (2003, Section 4.7) and Kuhn
19 481 and Johnson (2018, Section 4.8) provide a useful discussion.

30 482 **5. Examples of current fitting procedures**

31 483 In this section we make use of three different data sets to illustrate the strengths and weaknesses
32 484 of fitting limit lines using some of the simpler methods introduced above. For conciseness, we
33 485 have focussed on those simpler methods. The issues that arise using simpler methods also apply
34 486 to, and would inform the application of, more advanced statistical procedures. The application
35 487 here of simpler methods does not imply that more sophisticated approaches could not be explored
36 488 beneficially in the case of these examples.

37 489 The first example consists of a complex of several data sets which, taken together, define a visual
38 490 upper limit line for which an upper limit is expected from theory. This example is used to
39 491 demonstrate the use of three relatively simpler methods together with fitting of a theoretical function
40 492 that makes use of the empirical data.

41 493 The second example consists of a single data sets that is inadequate to clearly define a visual
42 494 upper limit line, although a limit is reasonably expected from prior studies. This example is used to
43 495 demonstrate the use of three relatively statistically robust methods.

44 496 The third example consists of a single data set for which the variance in y increases rapidly as the
45 497 value of x increases, and both upper and lower limit lines are required. Solutions derived using a
46 498 simple robust method are contrasted to inspection functions.

47
48
49
50
51
52
53
54
55
56
57
58
59

60 500 *Example 1: Catastrophic outburst floods from dammed lakes*

1
2
3 501 Figure 2 serves as an example of the issues that arise from fitting limit lines using Inspection,
4 502 Selective Regression and application of a data-informed Theoretical Limit. The data sets
5 503 collectively represent the relationship between measured volumes of released lake water (V) and
6 504 the estimates of the peak discharge (Q_p) downstream due to catastrophic lake failure (O'Connor
7 505 *et al.*, 2013). It may be expected that variation in breaching mechanisms, channel geometry and
8 506 roughness (amongst other controls) will mediate the downstream translation of the flood wave so
9 507 that different peak discharge values might be obtained for the same initial lake volumes. However,
10 508 if the discharge from the lake is constrained by the initial geometry of the

15 509

16 510 **Figure 2 near here**

17
18
19
20 511 eroding breach (e.g. critical flow control), or by the way the flood translates down system, there
21 512 should be an upper limit to the scatter of peak discharge values. The data in Fig. 2 considered
22 513 collectively, or as separate data sets, provide some support for the critical flow control as is detailed
23 514 below.

24 515

25
26
27 516 *Inspection and selective regression:* The green curve is fitted by inspection ($Q_p = 0.1286V^{0.83}$) to
28 517 pass through the three outlying 'landslide group' data points. The four green points attract attention
29 518 because, on log-log coordinates, the four points trace out a straight line lying above the main data
30 519 spread. Having fitted the green curve, the red 'constructed group' data point (lying above the green
31 520 curve) is defined as an outlier *a posteriori*, by the simple fact that it lies above the green limit line.
32 521 Note that this definition of the outlier is unsatisfactory given that robust methods (noted above) are
33 522 available to determine leverage. All other data points are included within the limit, but forward
34 523 extrapolation of the limit line means that the curve increasingly deviates away from the observed
35 524 data. A curve fitted through the four 'landslide group' data points (selective regression on a data
36 525 sub-set), using least-squares regression, provides a similar curve ($Q_p = 0.1168V^{0.83}$) and is
37 526 preferred to the eye-fitted curve for reasons explained prior.

38 527

39
40
41
42
43
44
45
46 528 *Selective regression with optimal a^* :* A least-squares regression of the 'physical model' data
47 529 defines the trend of that data set which, when extrapolated forwards (not shown) passes through
48 530 the centre of the mass of other data sets. This consilience between the two groups of data suggest
49 531 that the small-scale model results reproduce well the central tendency of behaviour of large natural
50 532 dam failures across several orders of magnitude. Interestingly, such an extrapolation might define
51 533 an upper limit line for 'Ice dams – subglacial tunnelling', although we do not explore the implications
52 534 herein. However, to define a limit line for the majority of data, the trend of the 'physical model' data
53 535 can be adjusted by adding increments to the intercept value, a^* , until sufficient data points fall
54 536 below the limit. In the example provided, the intercept value is increased (*Selective regression*
55 537 *with optimal a^**) by a factor of ten such that although ten data points lie outside the limit, the black

line provides a reasonably satisfactory visible limit to the data spread, notably that of the 'constructed group' and 'moraine group' data sets. A small increase in the intercept value would readily include seven more data points leaving only three as outliers. By adjusting the intercept value the exponent of the trend line is preserved, implying that the central tendency growth function for 'physical model' data also can define the behaviour of data at the upper limit to the larger-scale dam-break data. By such systematic exploration of central tendency and limits, consideration can be given to: i) the relationship of one data set to another; and ii) the consistency of data point plotting positions within the individual data sets. Further, iii) the positions of some individual data points come under scrutiny and; iv) possible theoretical constraints on the data plotting positions may become evident.

Theoretical Limit: A theoretical critical flow control might be considered to provide an upper limit to the data spread in Fig. 2. The theoretical derivation is provided as Supplement 1 but the basic facts are as follows. Failure of earthen and ice dams often is associated with initial establishment of a critical-flow depth (h_c) at the breach that determines the peak outflow discharge (Walder and O'Connor, 1997; O'Connor and Beebee, 2009). Larger volume (V) lakes tend to have greater depths (h) and so have the propensity to develop rapid failures with greater critical flow depths; thus $h_c \propto h$. Assuming that the outflow breach, and thus the critical flow depth, will be larger for larger water bodies, the maximum discharge Q_p should be proportional to the lake volume efflux (V). O'Connor and Beebee (2009) showed that a critical flow control can be approximated as:

$$Q_p = pg^{1/2}h_c^{5/2}, \quad (1)$$

where g is the acceleration due to gravity and p is a proportionality coefficient.

The V -data in Fig. 2 are used to calculate Q_p , so defining the values of h_c and p in Equation 1 can be seen as a fitting procedure, rather than just adding a theoretical function to the graph. Equation 1 provides a theoretical basis for the slope (b) of the green limit line. The slope of Eq. 1 is practically coincident with both the line drawn by inspection and the least-squares power function ($b = 0.83$) obtained using selective regression on a data sub-set (as reported above), matching the position of these two curves when $p = 1.0$. It is beyond the scope of this paper to discuss the reasons why the limit line constructed using theory has a steeper gradient than that devised using *selective regression with optimal a^** . Nevertheless, fitting the various limit lines leads to considerations such as that the theory applied may be too restrictive, or the small-scale physical model data may not adequately represent larger natural systems.

Example 2: Lichen growth curve to date flood deposits

1
2
3 574 Figure 3 serves as an example of the issues that arise from fitting limit lines using parametric
4 575 mixture modelling. The data considered (Carling, 1987) define the relationship between the
5 576 diameter of the largest lichen thalli on dated gravestones in Teesdale, northern England (Fig. 1B).
6 577 Such lichen growth curves can be used to date the surface of rocks that have been transported by
7 578 floods or glaciers in the same region for which the calibration data were obtained. The supposition
8 579 is that geophysical flows transport, abrade and destroy any pre-existing lichens, such that lichen
9 580 growth only occurs once the rocks are stable in a deposit. In this manner flood gravel bars ~~and~~
10 581 ~~glacial moraines~~ can be dated. The species of lichen (*Huillia albocaerulescens*) used by Carling
11 582 (1987) tends to produce circular thalli which, after an initial rapid growth phase of a few years only,
12 583 tend to steadily increase linearly in size with age. Eventually lichens reach senescence, at which
13 584 time lichen thalli cease to grow, grow more slowly, or being to break-up. Consequently, any
14 585 maximum linear growth function can only be extended to a given x:y breakpoint value beyond
15 586 which maximum growth does not apply (Cooley *et al.* 2006). Beyond this point, either a separate
16 587 lower-gradient function is fitted for the senescence phase, or, if a single function is fitted it must
17 588 account for the growth and senescence phases (Innes, 1983). In ideal growth conditions, lichens
18 589 will achieve a maximum diameter during the rapid growth phase. Data scatter occurs below an
19 590 expected upper limit to the x:y data pairs occurs for a number of reasons, including: pollution, the
20 591 date on the gravestone being added some time after erection; differences in the rock type, aspect,
21 592 and occasional cleaning of gravestones.

22 593
23 594 *Box and Lucas:* The data shown in Fig. 3 produces an upper limit line (blue curve) when using the
24 595 Originlab® procedure, that is of the same form as a conventional least-squares exponential fit
25 596 (orange curve) through all the data. Both curves are constrained to have an origin at T equals
26 597 zero, although other intercepts could be specified. A linear least-squares zero-intercept fit to all
27 598 the $T \leq 190$ data pairs (not shown), to represent only the growth phase, statistically would be a
28 599 less good fit ($r^2 = 0.31$) than the orange curve. The fitted limit is that which maximizes the r^2 value
29 600 for eight outer points, so other curves could be selected if desired. The points that lie just above
30 601 the D_{max} exponential solution were determined to do so by the final choice of the curved fitted. The
31 602 fitted line intuitively is acceptable as it encloses 93% of the data points, but a higher curve could
32 603 equally be obtained to enclose more data points.

33 604
34 605 *Mixture modelling:* An expectation–maximization (EM) algorithm was used to fit the red curve in
35 606 Fig. 3 following the mixture model of Maller *et al.*, (1983). The least-squares trimming method of
36 607 Maller *et al.*, (1983) leads to a solution (green curve) that is similar to selective regression (which
37 608 fits a least-squares function to an arbitrary selection of data points), but the degree of objectivity in
38 609 curve fitting is greater using mixture modelling. The solution is not uniquely determined, but the
39 610 accepted fitted line usually is taken to be the solution that includes the greatest number of data
40 611 points. In the case of the data in Fig. 3, a limit was derived after eight iterations which enclosed

1
2
3 612 95% of all data points and which passes through a further 4% leaving two points just above the
4 613 curve. The fitted curve: $D_{max} = 0.815T + 24.33$ lies slightly below the red curve fitted using EM
5 614 algorithm which enclosed all data points.
6
7

8 615

9 616 **Figure 3 near here**

10 617

11
12 618 As lichens often exhibit initial rapid growth, followed by a linear growth phase, followed by an
13 619 exponential decline during senescence (Cooley *et al.* 2006), a bipartite or tripartite limit line might
14 620 be preferable, although in the case of the data in Fig. 3 there are inadequate data to define a
15 621 separate senescence phase. However, it would be more satisfying if recourse was made to
16 622 biologically-based theoretical models of lichen growth (Childress and Keller, 1980) to determine
17 623 what form of function should be fitted that mimics the growth of lichens.
18
19
20 624

21 625

22 626 *Example 3: Variation in energy expenditure required to fracture pebbles*

23
24 627 Figure 4 serves as an example of the issues that arise from fitting limit lines using Inspection and
25 628 and Iterative Selective Regression. Figure 4 reproduces the data shown in Fig. 1A, with additional
26 629 limit lines fitted. The data published originally by Tuitz *et al.* (2012) were presented in this graphical
27 630 context by Carling & Fan (2020) The data represent the variation in experimentally-derived energy
28 631 expenditures recorded using a laboratory point-load test to fracture river pebbles. It is known from
29 632 theory and empirical measurements in prior published studies of fracture processes that the energy
30 633 should increase in a linear manner for the range of pebble sizes considered here. However, as
31 634 pebble size increases the number and complexity of flaws in the pebbles also increases such that
32 635 the variance in the y-data increases as a function of x. Carling & Fan (2020) only wished to draw
33 636 attention to the data spread and eye-fitted the red-dotted lines to delimit the data spread. The
34 637 lower and upper blue fitted limit lines were obtained after seven and nine iterations respectively
35 638 using iterative selective regression.
36
37
38
39
40
41
42

43 639 **Figure 4 near here**

44 640

45 641 **6. Concluding Discussion**

46 642

47 643 Researchers sometimes wish to define boundaries, upper or lower limits to samples of data, and
48 644 hence to the distributions from which those samples are drawn. In choosing an approach to
49 645 achieve this, the researcher should be as specific as possible about the objective of their data
50 646 analysis. Consideration should be given as to how the inferences derived from the analysis will be
51 647 used further to inform decisions. In some fields, including hydrology and environmental
52 648 engineering, there are specific concerns regarding characterisation of extreme values of the data-
53 649 generating process. In these areas, techniques motivated by extreme value theory are relatively
54 650 commonplace to quantify the (joint) tails of distributions from samples, and to estimate extreme
55
56
57
58
59
60 651 quantiles including upper bounds for conditional distributions such as $Y|x$. However, in many other

1
2
3 651 fields, estimation of boundaries or limit lines has received little or no attention. Rather weak *ad-*
4 652 *hoc* methods, making limited use of available data and quantitative modelling, have been applied.
5
6 653 On occasion, statistical methods such as linear regression, devised to characterize the general
7
8 654 nature of data spread $Y|x$ have been adapted to locate possible limit lines. Rarely have statistical
9
10 655 approaches which specifically seek to characterise the tail $Y|Y > u, x$ been used. Often a limited
11
12 656 number of observations precludes statistical modelling. Specifically, for the applications illustrated
13
14 657 in Figures 1(B), 2, 3 and 4, sample size is sufficient to attempt relatively simple regression models
15
16 658 for $Y|x$, including quantile regression; however, it would not be feasible to quantify the conditional
17
18 659 tail $Y|Y > u, x$ using extreme value analysis. Sometimes, weaker *ad-hoc* methods are adopted
19
20 660 because of a lack of awareness or appreciation that more principled approaches may be useful.
21
22 661 In general, *ad hoc* methods should not be used in cases where more principled statistical
23
24 662 procedures can be applied, because the latter are clearly defined mathematical models making
25
26 663 use of available data, are reproducible and allow quantification of uncertainty. Whereas *ad hoc*
27
28 664 methods introduce uncertainty with respect to interpretation, adoption of statistical procedures
29
30 665 allows both authors of articles and readers to further explore the implications of the fitted functions
31
32 666 in a rational manner.

667

33
34 668 In the absence of theoretical knowledge as to the form of a limit line, the qualitative procedure of
35
36 669 inspection is a useful initial means to consider the likely form of a function. Indeed, the intuitive
37
38 670 understanding of how the data behaves can assist in statistical model formulation, yet at the same
39
40 671 time inspection can lead to false inferences as to the likely behaviour of a limit. The quantitative
41
42 672 nature of data allows objective fitting of a statistical function, which can then be compared with the
43
44 673 intuitive expectations of the analyst. Given that a variety of statistical models are available, it is
45
46 674 important to consider at the outset the purpose of the fitting exercise and to choose the method
47
48 675 that is most appropriate to satisfy the objective. Fitting statistically derived limit lines is especially
49
50 676 powerful in those cases where the theoretical limit is either well-known or the behaviour is
51
52 677 reasonably expected. In these cases, the close agreement of the statistically fitted limit with a
53
54 678 theoretically derived line can be confirmatory. In contrast, significant discrepancies between the
55
56 679 two curves may indicate deficiencies with the data sample: additional data may be required, or the
57
58 680 quality of existing data may be suspect. Discrepancies may also highlight theoretical or model
59
60 681 inadequacies: the possibility that other covariates are affecting y - or x -values, or that the theory
62
682 may need revision.

683

684 In the examples provided herein (section 5) it is evident that the application of different methods
685 produces different limit lines. In some applications these discrepancies may not be significant. As
686 previously noted, the identification of extreme behaviour within environmental systems can be very
687 important for instance in hazard mitigation. In such critical situations the development of limit lines
688 rationally informed by empirical evidence, statistical and physical theory is preferable. Although

1
2
3 689 this conclusion may seem obvious, there are many examples in the literature of limit lines fitted
4
5 690 without consideration of existing theory. For example, surprisingly, limit lines are often fitted to
6
7 691 define the relationship between the maximum flood discharges generated from given catchment
8
9 692 areas without consideration of the maximum probable flood (MPF). The MPF is the theoretical
10
11 693 expectation (e.g. Shalaby, 1994; USFERC, 2001) and it would be informative to compare the
12
13 694 statistically derived flood limit lines with the theoretical functions. Where theory is unavailable,
14
15 695 consideration should be given as to whether the application of different methods tends to lead to
16
17 696 convergence in terms of the form and trend of several limit lines. In general however, identifying
18
19 697 the subset of methods that provide consistent estimates of limit lines is likely only to be possible
20
21 698 once the details of the problem and data have been understood. Building an appreciation for the
22
23 699 relative performance of different methodologies via simulation study for a specific problem type is
24
25 700 useful and standard practice in the statistics literature. However, the number of potential problem
26
27 701 types is huge, and therefore the specifics of the problem of interest first need to be clearly defined
28
29 702 before the simulation study is undertaken.

703

26 704 The use of advanced statistically methods in contrast to simple ones readily can be justified
27
28 705 (Jomelli *et al.* 2010), especially when there is plentiful empirical evidence. Not least, given the
29
30 706 inevitable ambiguity in fitting of limit lines, it is important to reason systematically whilst recognizing
31
32 707 the uncertain evidence that even large data sets offer (e.g. using Bayesian analysis). However,
33
34 708 situations occur where the x-y data points are few, or their disposition on the scatter plot render
35
36 709 the application of sophisticated methods impracticable or impossible. Such situations usually
37
38 710 indicate that additional data are required, or that stronger assumptions about the data-generating
39
40 711 process are necessary. Regardless, the procedure used to fit a limit line should be documented
41
42 712 sufficiently clearly that limit line estimation given a sample of data can be reproduced with
43
44 713 confidence. Fitting a limit by Inspection alone rarely can be justified.

714

42 715 The advantage of a statistical approach in general is that it provides a rational, reproducible basis
43
44 716 for inference, and hence a sound basis for learning: different practitioners working independently
45
46 717 can be reasonably expected to make the same inference given a sample of data. The performance
47
48 718 of a model is dependent on the quality of information used to infer it. It is not reasonable in general
49
50 719 to expect that a statistical model provides a "better result" than a visual fit, since a well-informed
51
52 720 visual fit may be superior to a badly specified statistical model. However, it is also self-evident than
53
54 721 an ill-informed visual fit can lead to spectacularly bad inferences.

722

54 723 The outline taxonomy or road map provided in Section 2 provides an overview of the range of
55
56 724 statistical methodologies available for estimation of limit lines, and references to statistical texts
57
58 725 which explain methodologies in more detail. Choice of the appropriate methodology will be problem
59
60 726 specific. When dealing with an unfamiliar problem, seeking the advice of a statistician is likely to

1
2
3 727 **be beneficial.** Given the uncertainty that can pertain to model fitting, we conclude by providing
4 728 some signposts that may assist in the decision-making process of **limit line** fitting:

5 729

6 730

- Define the objective of the analysis: for what purpose will the fitted limit line be used?

7 731 Consider how this informs the analysis to be undertaken

8 732

- Assess the data to hand, the characteristics of the measurement used to gather data, and likely sources of uncertainty. Are the measurements independent (given covariates)? Are the data representative? What is the potential for gathering further relevant data?

9 733

- Determine if theory allows the form of the limiting function to be defined

10 734

- Determine whether a statistical model can be adopted for the data-generating process and fitted to the data. Limit lines may then be estimated using the fitted statistical model. What form of statistical model is likely to be more appropriate? Otherwise consider what form of limit line curve might be appropriate from knowledge of the system behaviour

11 740

- Assess the appropriate level of sophistication of the statistical model or limit line curve, guided by parsimony. Is it likely that (unknown, unmeasured) covariates are in play? Should breakpoints be considered?

12 741

- In fitting the statistical model or limit line, always assess fitting performance using diagnostic plots and tools. Assess potential outliers.

13 742

- Seek to quantify uncertainties in the fitted model (line), and propagate those uncertainties to subsequent decisions made using the fitted model (line)

14 743

15 744 **7. Acknowledgements**

16 745 **XXXXXX** acknowledges the receipt of China Postdoctoral Science Foundation Grant No.

17 746 2020M670435. Software for the trimming method of Maller et al (1983) is provided at Carling *et al.* (2021), and for simple non-stationary extreme value analysis at Jonathan and Ewans (2021).

18 747 We are grateful to the Associate Editor, Karen Anderson, and two anonymous reviewers for their comments which substantially improved the presentation of the results.

19 748 **Declaration of conflicting interests**

20 749

21 750 The authors declared no potential conflicts of interest with respect to the research, authorship, and publication of this article.

22 751

23 752 **Funding**

24 753

25 754 The author(s) received no financial support for the research, authorship, and/or publication of this article.

26 755

27 756 **ORCID iD**

28 757

29 758

30 759

31 760 **References**

- 1
2
3 770
4 771 Aitkin M and Tunnicliffe Wilson GT (1980) Mixture models, outliers and the EM algorithm.
5 772 *Technometrics*, 22: 325–31.
6
7 773 Bagnold RA (1966) An approach to the sediment transport problem from general physics. *U.S.*
8 774 *Geological Survey Professional Paper*, 422-1, 37 pp.
9
10 775 Bagnold RA (1980) An empirical correlation of bedload transport rates in flumes and natural
11 776 rivers. *Proceedings of the Royal Society of London, A*, 372: 453-473.
12 777
13
14 778 Bishop C (2006) *Pattern Recognition and Machine Learning*. Springer.
15
16 779 Box GEP Wilson KB (1951) On the experimental attainment of optimum conditions. *Journal of the*
17 780 *Royal Statistical Society, Series B*, 13: 1-38.
18
19 781 Box GEP Lucas HL (1959) Design of experiments in non-linear situations. *Biometrika* 46: 77-80.
20
21 782 Brereton RG. (2009) *Chemometrics for Pattern Recognition*. John Wiley and Sons: Chichester.
22
23 783 Brereton RG Lloyd GR (2014) Partial least squares discriminant analysis: taking the magic away.
24 784 *Journal of Chemometrics* 28: 213–225.
25
26 785 Cade BS (2017) Quantile regression applications in ecology and the environmental sciences.
27 786 Pages 429-454 in Koenker R *et al.* eds. *Handbooks of Modern Statistical Methods:*
28 787 *Handbook of Quantile Regression*. Chapman & Hall/CRC.
29
30 788 Carling, PA (1987) Lichenometric dating applied to flood deposits. In: Beschta RL Blinn T, Grant GE
31 789 G.G.Ice and Swanson FJ (eds) *Proceedings of a Symposium on Erosion and Sedimentation*
32 790 *in the Pacific Rim, Corvallis, pp.395-396.*
33
34 791
35
36 792 Carling PA (1989) Bedload transport in two gravel-bedded streams. *Earth Surface Processes and*
37 793 *Landforms* 14: 27-39.
38
39 794 ~~Carling, PA (1987) Lichenometric dating applied to flood deposits. In: Beschta RL Blinn T, Grant GE~~
40 795 ~~G.G.Ice and Swanson FJ (eds) *Proceedings of a Symposium on Erosion and Sedimentation*~~
41 796 ~~*in the Pacific Rim, Corvallis, pp.395-396.*~~
42
43 797 Carling PA Fan W (2020) Particle comminution defines megaflood and superflood energetics.
44 798 *Earth-Science Reviews* 204: 103087
45
46 799 Carling PA Jonathan P Teng S (2021) Spreadsheet software for the trimming method of Maller et
47 800 al. (1983). <https://github.com/ygraigarw/LimitLines>
48
49 801 Castellarin A (2007) Probabilistic envelope curves for design flood estimation at ungauged sites.
50 802 *Water Resources Research* 43: W04406, doi:10.1029/2005WR004384
51
52 803 Chavez-Demoulin, V. and Davison, A.C. (2005) Generalized additive modelling of sample
53 804 extremes, *J. Roy. Statist. Soc. Series C: Applied Statistics* 54, 207-222.
54
55 805 Childress S and Keller JB (1980) Lichen growth. *Journal of Theoretical Biology*, 82: 157-165.
56
57 806 Coles, S (2001) *An Introduction to Statistical Modeling of Extreme Values*, Springer.
58
59 807 Cook RD Weisberg S (1982). *Residuals and influence in regression*. Chapman and Hall.
60

- 1
2
3 808 Cooley D Naveau P Jomelli V Rabatel A and Grancher D (2006) A Bayesian hierarchical extreme
4 809 value model for lichenometry. *Environmetrics* 17: 555–574.
5
6 810 Davison AC (2003) *Statistical Models*. Cambridge University Press.
7
8 811 Davison AC and Ramesh NI (2000) Local likelihood smoothing of sample extremes. *J. R. Statist.*
9 812 *Soc. B*, 62: 191–208.
10
11 813 Davison AC Smith RL (1990) Models for exceedances over high thresholds, *J. R. Statist. Soc. B*.
12 814 52: 393–493.
13
14 815 Dixon SJ Brereton RG (2009) Comparison of performance of five common classifiers represented
15 816 as boundary methods: Euclidean distance to centroids, linear discriminant analysis,
16 817 quadratic discriminant analysis, learning vector quantization and support vector machines,
17 818 as dependent on data structure. *Chemometrics and Intelligent Laboratory Systems* 95: 1-
18 819 17.
20 820 Eberhardt LL Thomas JM (1991) Designing environmental field studies. *Ecological Monographs*
21 821 61: 53-73.
22
23 822 Gaume E Bain V Bernardara P Newinger O Barbuc M Bateman A Blaškovićová L Blöschl G Borga
24 823 M Dumitrescu A Daliakopoulos I Garcia J Irimescu A Kohnova S Koutroulis A Marchi L
25 824 Matreata S Medina V Preciso E Sempere-Torres D Stancalie G Szolgay J Tsanis I Velasco
26 825 D and Viglione A (2009) A compilation of data on European flash floods. *Journal of*
27 826 *Hydrology* 367: 70–78.
29 827 Gelman A Carlin JB Stern HS Dunson DB Vehtari A and Rubin DB (2013) *Bayesian Data Analysis*,
30 828 Third Edition. Chapman and Hall/CRC.
31
32 829 Good PI (2006) *Resampling Methods: a Practical Guide to Data Analysis*. Birkhauser.
33
34 830 Hall P Tajvidi N (2000) Nonparametric analysis of temporal trend when fitting parametric models
35 831 to extreme-value data. *Statist. Sci.*, 15: 153–167.
36
37 832 Hao L Naiman DQ (2007). *Quantile Regression*. London: Sage Publications.
38 833
39 834 Hesterberg TC (2015) What teachers should know about the bootstrap: Resampling in the
40 835 undergraduate statistics curriculum. *American Statistician* 69: 371–86.
41 836
42 837 Innes JL (1983) Development of lichenometric dating curves for Highland Scotland. *Transactions*
43 838 *of the Royal Society of Edinburgh: Earth Sciences* 74: 23-32.
44
45 839 Joe H (2014) *Dependence Modelling with Copulas*. CRC Press.
46
47
48 840 Jomelli J Naveau P Cooley D Grancher D Brunstein D and Rabatel A (2020) A response to
49 841 Bradwell's commentary on "Recent statistical studies in lichenometry". *Geografiska*
50 842 *Annaler: Series A, Physical Geography* 92:485-487.
51
52 843 Jonathan P and Ewans K (2021) MATLAB code for simple non-stationary extreme value
53 844 analysis, estimated using MCMC.
54 845 <https://github.com/ygraigarw/SimpleNonstationaryExtremesBayesian>.
55
56 846 Kaiser MS Speckman PL and Jones JR (1994) Statistical models for limiting nutrient relations in
57 847 inland waters. *Journal of the American Statistical Association* 89: 410-423.
58
59 848 Koenker R (2005). *Quantile Regression*. New York: Cambridge University Press.
60

- 1
2
3 849 Kuhn M and Johnson K (2018). Applied predictive modeling. Springer.
4
5 850 Lehr D and P Ohm P (2017) Playing with the data: what legal scholars should learn about machine
6 851 learning. U C Davis Law Review 51: 653-717.
7
8 852 Maller RA de Boer ES Joll LM Anderson DA and Hinde JP (1983) Determination of the maximum
9 853 foregut volume of Western Rock Lobsters (*Panulirus cygnus*) from field data. Biometrics
10 854 29: 543-551.
11 855 McLachlan GJ Lee SX and Rathnayake SI (2009) Finite mixture models. Annual Review of
12 856 Statistics and Its Application 6: 355-378.
13
14 857 Molinaro AM Simon R and Pfeiffer RM (2005) Prediction error estimation: a comparison of
15 858 resampling methods. Bioinformatics 21: 3301-7.
16 859 O'Connor JE and Beebee RA (2009) Floods from natural rock-material dams. In: Burr DM. Carling
17 860 PA Baker VR (eds) Megaflooding on Earth and Mars. Cambridge University Press,
18 861 Cambridge, UK, pp. 128-171.
19
20 862 O'Connor JE Clague JJ Walder JS Manville V and Beebee RA (2013) Outburst Floods. In: Shroder
21 863 JF (Editor-in-Chief), Wohl E (Volume Editor). Treatise on Geomorphology, Vol 9, Fluvial
22 864 Geomorphology, San Diego: Academic Press. pp. 475-510.
23
24 865 Pawitan, Y. (2001) In All Likelihood: Statistical Modelling and Inference Using Likelihood. Oxford.
25
26 866 Ramesh NI Davison AC (2002) Local models for exploratory analysis of hydrological extremes.
27 867 Journal of Hydrology 256: 106-119.
28
29 868 Reistad M Breivik O Haakenstad H Aarnes OJ Furevik BR and Bidlot JR (2011). A high-resolution
30 869 hindcast of wind and waves for the North Sea, the Norwegian Sea, and the Barents Sea.
31 870 Journal of Geophysical Research 116: 1-18.
32
33 871 Ryan SE Porth LS Troendle CA (2002) Defining phases of bedload transport using piecewise
34 872 regression. Earth Surface Processes and Landforms 27: 971-990.
35
36 873 Shalaby AI (1994) Estimating probable maximum flood probabilities. Journal of the American
37 874 Water Resources Association 30: 307-318.
38
39 875 Shirazi M Khademalrasoul A Ardebili SMS (2020) Multi-objective optimization of soil erosion
40 876 parameters using response surface method (RSM) in the Emamzadeh watershed. Acta
41 877 Geophysica 68: 505-517.
42
43 878 Tarolli P Borga M Morin E and Delrieu G (2012) Analysis of flash flood regimes in the North-
44 879 Western and South-Eastern Mediterranean regions. Nat. Hazards Earth Syst.Sci. 12:
45 880 1255-1265.
46
47 881 Tuitz C Exner U Frehner M and Grasemann B (2012) The impact of ellipsoidal particle shape on
48 882 pebble breakage in gravel. International Journal of Rock Mechanics & Mining Sciences 54:
49 883 70-79.
50
51 884 USFERC (2001) United States Federal Energy Regulatory Commission, 2001. Determination of
52 885 the probable maximum flood (Chap. VIII). In *Engineering Guidelines for the Evaluation*
53 886 *of Hydropower Projects*. Washington (DC): United States Department of Energy, pp 121.
54
55 887 Walder JS and O'Connor JE (1997) Methods for predicting peak discharge of floods caused by
56 888 failure of natural and constructed earthen dams. Water Resource Research 33: 2337-2348.
57
58 889 Wetherill GB Duncombe P Kenward M Kollerstrom J Paul SR and Vowden BJ (1986) Regression
59 890 analysis with applications. Springer.
60

891 Yu C Chen K and Yao W (2015) Outlier detection and robust mixture modeling using nonconvex
892 penalized likelihood, *Journal of Statistical Planning and Inference* 164: 27-38.

893 Zanini E Eastoe E Jones M Randell D and Jonathan P (2020) Covariate representations for non-
894 stationary extremes. *Environmetrics* e2624.

895

896

897 **Figure Captions**

898 *Figure 1: Examples of limit line fits. A) Central tendency in the relationship between the size to*
899 *pebbles and the energy required to break them is defined by least squares regression (blue curve).*
900 *Uncertainty in the energy required increases as a function of the pebble size. Limit lines (red) are*
901 *defined using Inspection (explained in text); B) Lichen growth curve: Central tendency defined by*
902 *zero-intercept (blue) regression curve; Limit line (grey) defined by simple linear regression with*
903 *adjusted intercept (explained in text) to enclose all data points. C) Significant wave height as a*
904 *function of wind speed at a location in the north-east Atlantic, with piecewise-linear quantile*
905 *regressions at the 0.9 quantile level fitted independently to the data below and above the median*
906 *x-value of 10 ms⁻¹; Pebble data from Tuitz et al., (2012); Lichen data from Carling (1987); wave*
907 *data from Reistad et al., (2011).*

908

909 *Figure 2: Empirical data define the relationship between the flood volume and the peak discharge*
910 *of water released from catastrophic failures of dammed lakes. Brown curve is the least-squares*
911 *fit to the physical model data; Limit line (black) fitted to all the data using selective regression with*
912 *optimal a*; Limit line (green) fitted by inspection of a data sub-set. The equivalent theoretical*
913 *equation, $Q_p = g^{1/2}(h_c)^{5.2}$, is essentially the same as the green line (see main text).*

914 *Figure 3: Empirical relationship between the date on gravestones and the diameter of lichen thalli*
915 *in 1986. Data from Carling (1986). The red curve was fitted using an EM algorithm. The green*
916 *curve was fitted using the Maller et al. (1983) trimming method. The blue curve was fitted using*
917 *the Box & Lucas (1959) method. The orange curve was fitted to all the data using a least-squares*
918 *exponential fit.*

919

920 *Fig. 4: Variation in experimentally-derived energy expenditures recorded using point load test*
921 *applied to fracture water-worn pebbles. Red curves were fitted by visual inspection. Blue curves*
922 *were fitted using selective regression.*

923

924 *Figure S1: Representation of the spread of x-y data wherein the variance of $Y|X = x$ increases as*
925 *a function of x . The data shows the burned area of forest (on logarithm base 10 scale) against*

1
2
3 926 *background atmospheric temperature, taken from Cortez and Morais (2007), available at*
4 927 <https://archive.ics.uci.edu/ml/datasets/Forest+Fires> .The red, orange and blue lines represent
5 928 *estimated linear quantile regression lines for the 0.9, 0.5 and 0.1 quantile levels. The black curves*
6 929 *illustrate Gaussian density fits to the conditional distribution of $Y|X = x$ for different choices of x .*
7
8
9

10 930

11 931

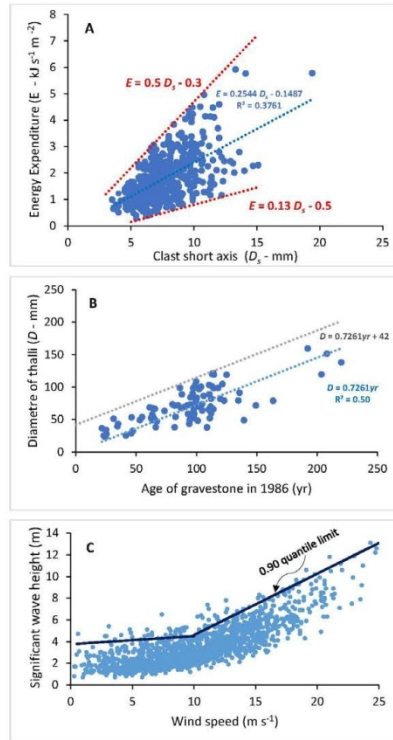
12 932

13 933

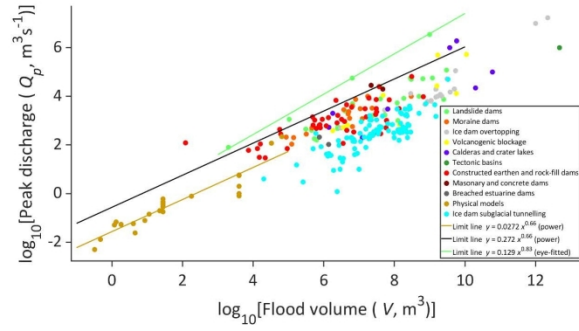
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

For Peer Review

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

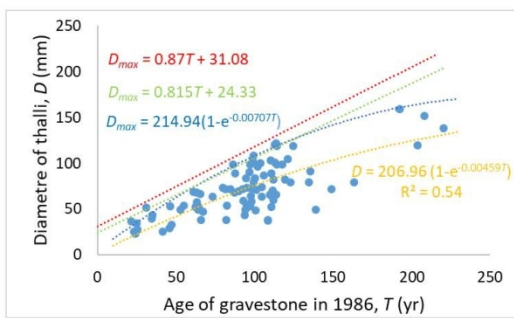


210x297mm (150 x 150 DPI)

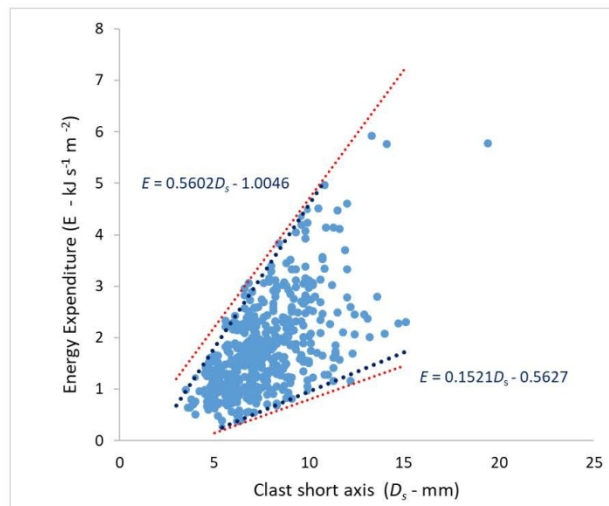


210x297mm (220 x 220 DPI)

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

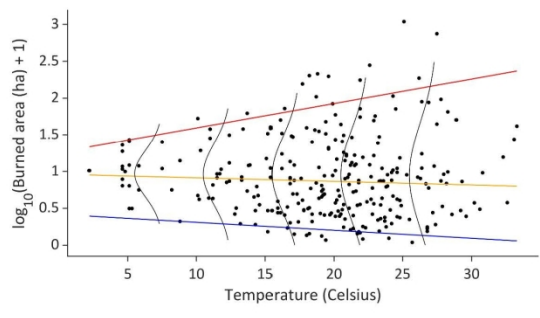


210x297mm (150 x 150 DPI)



210x297mm (150 x 150 DPI)

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60



210x297mm (221 x 220 DPI)

Supplement 1: Theoretical limit

A theoretical critical flow control might be considered as follows. Failure of earthen and ice dams often is associated with initial establishment of a critical-flow depth (h_c) at the breach that determines the peak outflow discharge (Walder and O'Connor, 1997; O'Connor and Beebee, 2009): Larger volume lakes tend to have greater depths (h) and so have the propensity to develop rapid failures with greater critical flow depths; thus $h_c \propto h$. Assuming that the outflow breach, and thus the critical flow depth, will be larger for larger water bodies, the maximum discharge Q_p should be proportional to the lake volume efflux (V). The scale values of h or h_c are unknown for the cases considered and have to be estimated, as follows. To derive estimates of h we note that many dammed lakes in V-shaped valleys have a geometry that approximates a tetrahedron, being deep at the dam, shallowing and narrowing up valley. An equilateral triangle with length, W , forms a vertical plane at the dam, with the apex of the tetrahedron denoting the upstream extent (L) of the lake, which has a triangle water-surface area. Assuming a regular tetrahedron ($W = L$), then $W = \sqrt{2}(3V)^{\frac{1}{3}}$ and, from Pythagorus' theorem the water depth at the dam face, $h = \sqrt{W^2 - \left(\frac{1}{2}W\right)^2}$.

For the landslide dams plotted in Fig. 2, O'Connor and Beebee (2009) determined:

$$h_c > 0.2V^{0.14} h^{0.58} \quad (1)$$

O'Connor and Beebee (2009) showed that a critical flow control can be approximated as:

$$Q_p = pg^{1/2}h_c^{5/2}, \quad (2)$$

where g is the acceleration due to gravity and p is a proportionality coefficient.

Given that h_c calculated from Eq. 1 represents a minimum value for landslide dams, the proportionality coefficient in Eq. 2 subsumes the fact that the true value of h_c likely will be greater than calculated using equation 1, but also reflects the fact that $L > W$ and $h_c \leq h$. Using values of h_c derived from the landslide dam V-data, Equation 2 provides the theoretical basis for the slope (b) of the green limit line. The slope of Eq. 2 is practically coincident with both the eye-drawn line and least-squares power function ($b = 0.83$) reported in the main text, matching the position of these two curves when $p = 1.0$. It is beyond the scope of this paper to discuss the reasons why the limit line constructed using theory has a steeper gradient than that devised using regression analysis. Nevertheless, fitting the various limit lines leads to considerations such as that the theory

1
2
3 38 applied may be too restrictive, or the small-scale physical model data may not adequately
4 39 represent larger natural systems.

7 40 **Supplement 2: Quantile regression**

9 41 **Figure S1 near here**

11 42
12 43 Figure S1 presents an example of regression data for which the variance of the response increases
14 44 as the predictor increases and there is a visual upper limit to the data spread. In this case, a
15 45 quantile regression model provides a suitable mechanism to estimate limit lines as discussed
17 46 below. The data represent burned area of forest as a function of background atmospheric
18 47 temperature.

20 48
21 49 In contrast, a simple linear regression model will not provide a suitable basis for estimation of limit
22 50 lines, essentially because the model is not appropriate to characterise the data-generating process.
23 51 For the data in Figure S1, a least squares regression for a response Y onto the predictor x models
24 52 the conditional mean $E(Y|X=x)$ as a linear function of x , and assumes that the distribution of
25 53 $Y|X=x$ is Normal with constant variance, not influenced by the value x . It does not therefore
26 54 capture the increasing conditional variance $var(Y|X=x)$ and more generally the conditional
27 55 distribution $Y|X=x$ of Y given x .

28 56
29 57 The black curves in Figure S1 represent the conditional densities of $Y|X=x$ for five specific values
30 58 of x . A set of densities for a comprehensive grid of values of x would provide a complete picture of
31 59 the conditional distribution of $Y|X=x$. Note that the conditional densities illustrated are assumed
32 60 Normal only for the purpose of illustration.

33 61
34 62 Figure S1 also shows fitted linear quantile regression models for quantile non-exceedance
35 63 probabilities 0.9, 0.5, and 0.10 (equivalently, the 90th, 50th, and 10th percentiles of the conditional
36 64 distribution $Y|X=x$ as a function of x). We might select the 0.1 and 0.9 quantile lines (or even more
37 65 extreme quantiles as appropriate) as limit lines.

38 66 Koenker (2005) and Hao and Naiman (2007) provide excellent introductions to the theory and
39 67 applications of quantile regression. Quantile regression software is available e.g. in PYTHON,
40 68 MATLAB and R.

41 69 We note that statistical models admitting heteroscedasticity would also provide appropriate
42 70 descriptions of data such as those in Figure S1, and hence yield principled estimates of a limit line.

43 71 **Supplement 3: Non-stationary extreme value analysis**

44 72

73 The objective of extreme value analysis is to estimate the tail of the distribution of a random
 74 variable. Theory (e.g. Beirlant *et al.* 2004) shows that independent occurrences of peaks over
 75 threshold for some random variable Y with a stationary distribution (satisfying a max-stability
 76 condition) are asymptotically generalised Pareto distributed with cumulative distribution
 77 function:

$$78 \quad F_{Y|Y>u}(y) = 1 - \left(1 + \frac{\xi}{\sigma}(y - u)\right)^{-1/\xi}$$

79 for shape $\xi \in (-\infty, \infty) \setminus 0$, scale $\sigma > 0$ and threshold u . When $\xi = 0$ the distribution takes the
 80 form $1 - \exp(-(y - u)/\sigma)$. If the distribution of Y is stationary only when conditioned on
 81 covariate X , we might choose to adopt a high quantile of the non-stationary conditional
 82 distribution $Y|Y > u(x), X = x$ for high threshold $u(x)$, as a function of x , as a limit line. In this
 83 case the conditional density for $Y|Y > u(x), X = x$ becomes

$$84 \quad F_{Y|Y>u(x),X}(y|x) = 1 - \left(1 + \frac{\xi(x)}{\sigma(x)}(y - u(x))\right)^{-1/\xi(x)}$$

85
 86 for $\xi(x) \neq 0$, and $1 - \exp(-(y - u(x))/\sigma(x))$ when $\xi(x) = 0$, where ξ , σ and u are now all
 87 functions of x .

88 When the value of ξ is negative, the distribution of Y has a finite upper limit. Thus, if there is
 89 specific prior knowledge that a finite upper limit exists, it might be appropriate to restrict
 90 estimates for the value of ξ to be negative.

91 In simple cases, it might be appropriate to adopt linear forms for these parameters, such as
 92 $\xi(x) = a_\xi + b_\xi x$, with similar descriptions for σ and u ; in general, more sophisticated
 93 parameterisations are needed, e.g. as described in Zanini *et al.* (2020).

94 Estimating a limit line using extreme value theory therefore requires the following procedure:
 95 (a) estimate an extreme value threshold $u(x)$ e.g. using an empirical quantile, or quantile
 96 regression, corresponding to some high quantile non-exceedance probability τ at covariate
 97 value x ; (b) assume a generalised Pareto model for exceedances of $u(x)$, and estimate
 98 generalised Pareto parameters $\xi(x)$ and $\sigma(x)$; and (c) estimate a limit line as an extreme
 99 quantile $Q(x)$ of the fitted generalised Pareto corresponding to a large non-exceedance
 100 probability τ^* near unity. The value of $Q(x)$ is obtained by solving the equation above for y at
 101 each x , setting the left-hand side to τ^* . When interest lies in the lower tail $Y|Y < u(x), X = x$,
 102 changing the sign of Y transforms the problem into the upper tail case just discussed.

1
2
3 103 Software (Jonathan and Ewans, 2021) provides a simple algorithm to estimate non-stationary
4 104 extreme value threshold $u(x)$ using quantile regression, and generalised Pareto models for
5 105 threshold exceedances. Linear forms for $\xi(x)$, $\sigma(x)$ and $u(x)$ are assumed, and the estimation
6 106 is performed by Bayesian Markov chain Monte Carlo (MCMC) inference using adaptive MCMC
7 107 (Roberts and Rosenthal 2009). Using the fitted models, limit lines can be estimated as
8 108 described in step (c) of the previous paragraph. The key steps in an extreme value analysis
9 109 of peaks over threshold using MCMC are (a) specification of reasonable prior distributions for
10 110 the parameters of $\xi(x)$ and $\sigma(x)$, and (b) diagnosis that the estimate of $\xi(x)$ is relatively
11 111 insensitive to choice of threshold $u(x)$ (and hence τ ; see Coles 2001).

12
13
14 112 Extreme value analysis is used widely in environmental science and engineering. It is also
15 113 used e.g. in lichenometry. Cooley et al. (2006) illustrates the use of a Bayesian hierarchical
16 114 generalised extreme value model for lichenometry. Jomelli et al. (2020) provide a useful
17 115 motivation for extreme value analysis in lichenometry.

18
19
20
21
22
23
24
25 116

26
27 117
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

118 **References**

- 119 Beirlant J. Goegebeur Y Segers J Teugels J (2004) *Statistics of extremes: theory and*
120 *applications*. Wiley.
121
- 122 Coles, S (2001) *An Introduction to Statistical Modeling of Extreme Values*, Springer.
- 123 Cooley D Naveau P Jomelli V Rabatel A and Grancher D (2006) A Bayesian hierarchical
124 extreme value model for lichenometry. *Environmetrics* 17: 555–574.
- 125 Cortez P Morais A (2007). A Data Mining Approach to Predict Forest Fires using Meteorological
126 Data. In J. Neves, M. F. Santos and J. Machado Eds., *New Trends in Artificial*
127 *Intelligence, Proceedings of the 13th EPIA 2007 - Portuguese Conference on Artificial*
128 *Intelligence*, December, Guimarães, Portugal.
- 129 Hao L Naiman DQ (2007). *Quantile Regression*. London: Sage Publications.
130
- 131 Jonathan P and Ewans K (2021) MATLAB code for simple non-stationary extreme value
132 analysis, estimated using MCMC.
133 <https://github.com/ygraigarw/SimpleNonstationaryExtremesBayesian>.
- 134 Jomelli J Naveau P Cooley D Grancher D Brunstein D and Rabatel A (2020) A response to
135 Bradwell's commentary on "Recent statistical studies in lichenometry". *Geografiska*
136 *Annaler: Series A, Physical Geography* 92:485-487.
- 137 Koenker R (2005) *Quantile Regression*. New York: Cambridge University Press.
- 138 O'Connor JE and Beebee RA (2009) Floods from natural rock-material dams. In: Burr DM.
139 Carling PA Baker VR (eds) *Megaflooding on Earth and Mars*. Cambridge University
140 Press, Cambridge, UK, pp. 128–171.
- 141 Roberts GO Rosenthal JS (2009) Examples of Adaptive MCMC. *Journal of Computational &*
142 *Graphical Statistics* 18: 349–67.
- 143 Walder JS and O'Connor JE (1997) Methods for predicting peak discharge of floods caused
144 by failure of natural and constructed earthen dams. *Water Resource Research* 33:
145 2337–2348.
146