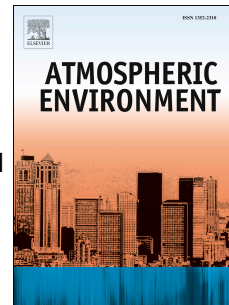


Journal Pre-proof



A strategy for modelling heavy-tailed greenhouse gases (GHG) using the generalised extreme value distribution: Are we overestimating GHG flux using the sample mean?

M.S. Dhanoa, A. Louro, L.M. Cardenas, A. Shepherd, R. Sanderson, S. Lopez, J. France

PII: S1352-2310(20)30237-5

DOI: <https://doi.org/10.1016/j.atmosenv.2020.117500>

Reference: AEA 117500

To appear in: *Atmospheric Environment*

Received Date: 30 November 2019

Revised Date: 21 March 2020

Accepted Date: 10 April 2020

Please cite this article as: Dhanoa, M.S., Louro, A., Cardenas, L.M., Shepherd, A., Sanderson, R., Lopez, S., France, J., A strategy for modelling heavy-tailed greenhouse gases (GHG) using the generalised extreme value distribution: Are we overestimating GHG flux using the sample mean?, *Atmospheric Environment* (2020), doi: <https://doi.org/10.1016/j.atmosenv.2020.117500>.

This is a PDF file of an article that has undergone enhancements after acceptance, such as the addition of a cover page and metadata, and formatting for readability, but it is not yet the definitive version of record. This version will undergo additional copyediting, typesetting and review before it is published in its final form, but we are providing this version to give early visibility of the article. Please note that, during the production process, errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.

© 2020 Published by Elsevier Ltd.

Dhanoa, M.S.: Conceptualization, Methodology, Analysis, Writing- original draft preparation. **Louro, A.:** Resources, Data curation, **Cardenas, L.M.:** Funding acquisition, writing mid-draft preparation, reviewing, editing. **Shepherd, A.:** Writing- final draft preparation, wrote response to peer review : **Sanderson, R:** Reviewing and editing. **Lopez, S.:** Reviewing, Editing, **France, J.:** Methodology, Reviewing, Editing.

Journal Pre-proof

1 **A strategy for modelling heavy-tailed greenhouse gases (GHG) using the generalised**
2 **extreme value distribution: Are we overestimating GHG flux using the sample mean?**

3

4 M.S. Dhanoa¹, A. Louro², L.M. Cardenas², A. Shepherd^{3*}, R. Sanderson⁴, S. Lopez^{5,6} and J.
5 France¹

6

7

8

9 ¹ Centre for Nutrition Modelling, Department of Animal Biosciences, University of Guelph,
10 Guelph ON, N1G 2W1, Canada.

11 ² Rothamsted Research, North Wyke, Okehampton, Devon, EX20 2SB, UK.

12 ³ Institute of Biological and Environmental Sciences, School of Biological Sciences, University
13 of Aberdeen, 23 St Machar Drive, Aberdeen AB24 3UU, UK

14 ⁴ Institute of Biological, Environmental and Rural Sciences, Aberystwyth University,
15 Gogerddan, Aberystwyth, Ceredigion, SY23 3EB, UK.

16 ⁵ Departamento de Producción Animal, Universidad de León, E-24007 León, Spain

17 ⁶ Instituto de Ganadería de Montaña, CSIC-Universidad de León, Finca Marzanas s/n,
18 Grulleros, 24346 Leon, Spain

19 *Corresponding author

20

21

22

23

24

25

26

27

28

29

30 **ABSTRACT**

31 In this study, we draw up a strategy for analysis of GHG field data. The distribution of
32 greenhouse gas (GHG) flux data generally exhibits excessive skewness and kurtosis. This
33 results in a heavy tailed distribution that is much longer than the tail of a log-normal
34 distribution or outlier induced skewness. The generalised extreme value (GEV) distribution is
35 well-suited to model such data. We evaluated GEV as a model for the analysis and a means
36 of extraction of a robust mean of carbon dioxide (CO₂) and nitrous oxide (N₂O) flux data
37 measured in an agricultural field. The option of transforming CO₂ flux data to the Box-Cox
38 scale in order to make the distribution normal, was also investigated. The results showed that
39 average CO₂ value estimates from GEV are less affected by data in the long tail compared to
40 the sample mean. The data for N₂O flux were much more complex than CO₂ flux data due to
41 the presence of negative fluxes. The estimate of the average value from GEV was much more
42 consistent with maximum data frequency position. The analysis of GEV, which considers the
43 effects of hot-spot-like observations, suggests that sample means and log-means may
44 overestimate GHG fluxes from agricultural fields. In this study, the arithmetic CO₂ sample
45 mean of 65.62 (mean log-scale 65.89) kg CO₂-C ha⁻¹ d⁻¹ was reduced to GEV mean of 60.14
46 kg CO₂-C ha⁻¹ d⁻¹. The arithmetic N₂O sample mean of 1.038 (mean log-scale 1.038) kg
47 N₂O-N ha⁻¹ d⁻¹ was reduced to GEV mean of 0.01571 kg N₂O-N ha⁻¹ d⁻¹. Our analysis
48 suggests that GHG data should be analysed using the GEV method, including a Box-Cox
49 transformation when negative data is present, rather than only calculating basic log and log-
50 normal summaries. Results of GHG studies may end up in national inventories. Thus, it is
51 necessary and important to follow all procedures that help minimise any biases in the data.
52

53 **Keywords:** nitrous oxide; carbon dioxide; generalised extreme value; Finney correction;
54 heavy-tailed data; skewness correction.

55 1. INTRODUCTION

56

57 Greenhouse gas (GHG) flux data from agricultural fields are difficult to measure precisely
58 because of their inherent spatial and temporal variability. This variability comes from
59 influencing factors such as soil moisture and underlying drainage, field aspect and slope, pH
60 and field distribution of dung or fertilizer. Hot-spots, or rather hot-moments (recorded peaks
61 are time peak rather than spatial peaks), in GHG data are a common occurrence and are a
62 cause of much nuisance for data analysis (Dixon et al. 2010; Loick et al. 2017). As a result,
63 data recorded on any time scale tend to include high and low peaks resulting in a skewed
64 distribution.

65 Although GHG emissions information can be extended by computer simulation using soil
66 biogeochemical cycling models, crucially the modelled data requires field data for calibration
67 and validation. Hence robust methods for analyses of field data are key to obtaining both
68 accurate field data and simulated data.

69 A common method of analysis is to transform skewed data to a log-scale. However, as
70 explained and illustrated in Dhanoa et al. (2016), skewness does not always mean a log-
71 normal distribution. Skewness caused by a few extreme values or outliers may be handled by
72 transforming data, e.g. using the Box and Cox (1964) system or the Finney (1941) correction.
73 If there are many outliers and the data transformation option fails, the generalised extreme
74 value (GEV) distribution offers an option. This is a very flexible distribution with only three
75 parameters to estimate, sometimes referred to as the Fisher–Tippett distribution after its
76 progenitors (Fisher and Tippett, 1928; Eastoe, 2017), though the common form used in several
77 versions of the GEV follows McFadden (1978).

78 The GEV is a class of probability distributions, incorporating a heavy-tailed distribution, that
79 we can fit to GHG data in order to extract metrics such as the mean and standard deviation.

80 Plant traits are generally positively skewed, and usually log-transformed. Edwards et al.
81 (2015) used GEV to determine the shape of seed mass distributions.

82 Küchenhoff and Thamerus (1996) used GEV in the extreme value analysis of Munich air
83 pollution data. Ercelebi and Toros (2009) also used GEV to model Istanbul air pollution (in
84 particular Ozone [O₃], Benzene [C₆H₆], nitric oxide [NO]). The interactions among these
85 affect N cycling, e.g. [NO + O₃ → N₂O + O₂].

86 Recently, for modelling air pollution data, Korkmaz (2015) described two-sided generalised
87 Gumbel (TSGG) distribution, which is a special case of GEV (type I distribution). Martins et
88 al. (2017) did extreme value modelling of air pollution data and compared results amongst
89 two large urban regions of South America. Beniston (2004) analysed the 2003 heat wave data
90 in Europe and showed an association with enhanced atmospheric GHG concentrations.

91 Battista et al (2016) used GEV to model urban concentrations of pollutants (legislated under
92 the Directive 2008/50/EC) in the city of Rome (Italy).

93 GEV is often applied to climatology. It is applied to changes in temperature and precipitation
94 extremes occurring as the effect of an increase in GHGs, to characterise event magnitudes
95 and frequencies (Kharin and Zwiers, 2004; Katz, 2010). Studies have so far tended to apply
96 GEV to the climate effects of GHG, rather than the sampled measurements of GHGs
97 themselves.

98

99 The purpose of this study is to assess the suitability of the GEV when analysing GHG data
100 from agricultural fields, which often contain larger than expected extreme values forming a
101 thick-tailed data distribution. Its purpose is also to show that the GEV method could
102 eliminate bias inherent in simple means of skewed GHG data, and to draw up a strategy for
103 analysis of GHG field data.

104

105 2. MATERIALS AND METHODS

106

107 2.1. Experimental design and data collection

108 The data set originated from a study conducted at Rothamsted Research, North Wyke, Devon,
109 UK (50:46:10N, 3:54:05W). The site is on a permanent grassland in a maritime climate
110 (mean annual temperature 9.6°C; mean annual precipitation 1056 mm).

111 Four treatments were tested: a) control with no nitrogen (N) fertilizer applications (CN); b)
112 digestate from anaerobic treatment of food waste (DG); c) ammonium nitrate (AN); d) cattle
113 slurry (SL) (Louro et al. 2013, Pezzola et al. 2012).

114 The soil is a silty clay loam, classed under the British soil classification as clayey typical non-
115 calcareous pelosol of the Halstow series and a stagni-vertic cambisol. The similar silty clay
116 loam used for the digestate sampling is a clayey non-calcareous Pelostagnogley of the
117 Hallsworth series.

118 The digestate (from Andigestion biogas plant in Holsworthy, UK) comprised food residues,
119 liquid waste from abattoirs and municipal waste from an anaerobic fermentation cycle lasting
120 50 d (or days).

121 Cattle slurry was collected from a dairy farm nearby the study site and spread at a rate of 80
122 kg total N ha⁻¹.

123 The applied ammonium nitrate comprised 34.5% N.

124 Chemical composition of the slurry and digestate can be seen in Table 1.

125 Further information on soil characteristics and chemical composition of the materials applied
126 can be found in Louro et al. (2013) and Pezzolla et al. (2012).

127 The four treatments were applied in a randomized block design with three replicate plots per
128 treatment on 8 September 2011, and applied to supply the equivalent rate of 80 kg N ha⁻¹ on
129 each occasion. Nitrous oxide (kg N₂O-N ha⁻¹ d⁻¹) and carbon dioxide (kg CO₂-C ha⁻¹ d⁻¹)

130 were measured in the 12 plots throughout 47 d between 12 September and 28 October 2011
131 using one dark non-transparent long-term chamber (LiCor 8100–104) per replicate plot
132 connected to a photoacoustic infrared gas monitor (Lumasense Technologies, INNOVA
133 model 1412i) and an infrared gas analyser (LI-COR Lincoln, Nebraska USA, model LI-
134 8100A). The flux was collected daily from the 12 chamber readings at 11:00H. There were
135 12 sets of data each with 47 observations.

136

137 *2.2. Analysis of GHG data with generalised extreme value (GEV) distribution*

138 A glossary of input parameters required for this study is listed in Table 2.

139 A kernel density plot (Sheather and Jones, 1991) for CO₂ and N₂O fluxes, showing the
140 position of observation frequency and the nature of skewness, is given in Figure 1. This
141 illustrates that the processes which cause GHGs to produce apparent outliers combine to give
142 data a heavy-tailed distribution (also known as thick-tailed, long-tailed, fat-tailed, etc.). When
143 one summarises such data, non-robust statistics such as the sample mean can be highly
144 inflated. The classic approach to deal with a skewed distribution is to check if it follows the
145 log-normal distribution. This is usually done by transforming the data to the log-scale and
146 then testing whether the transformed data follow the normal distribution. One complication
147 comes if the original data contains zeroes or negative values. In this case one must add a
148 positive constant equal to the sample minimum to make all data positive and one must also
149 add 1.0 where zero values are present. Once the constant has been applied, the data can be
150 transformed to the log-scale (see Dhanoa et al. 2016 for further information on log
151 transformation). It is worth noting here that the anti-logged value of the mean estimate from
152 the log-scaled data is not the same as the calculated arithmetic mean on the original scale,
153 rather the geometric mean. To calculate the mean on the original scale the Finney correction
154 must be applied. Finney (1941) showed that

$$155 \quad AM = e^{\hat{\mu} + \frac{\hat{\sigma}^2}{2}} = e^{\hat{\mu}} e^{\frac{\hat{\sigma}^2}{2}} \quad \text{eq. 1}$$

156 where AM is the arithmetic sample mean on the original scale and $\hat{\mu}$ and $\hat{\sigma}^2$ are the estimates
 157 of the sample mean and the variance on the log-scale, respectively. Any constant applied
 158 prior to logarithmic transformation should be subtracted.

159

160 An alternative option is to use a data transformation system such as the Box-Cox
 161 transformation (Box and Cox, 1964):

$$162 \quad y^* = \frac{(y+c)^{\lambda}-1}{\lambda} \quad \text{eq. 2}$$

163 Where y^* is the transformed data value, y is the actual data value, c is the positive constant
 164 added to make all data above zero and λ the transformation parameter, enabling the best
 165 approximation of a normal distribution curve. To perform this transformation, the value of λ
 166 must be estimated first. The algorithm to estimate λ [by numerical search] usually fails in the
 167 presence of negative values, so it is prudent to add a suitable constant as detailed above if
 168 negative or zero values are present. Having estimated the value of λ and checked if the
 169 transformed data follows a normal distribution, one now has the task of transforming the
 170 mean estimate back to the original scale. There are not many validated methods, with the
 171 exception of the method detailed by Taylor (1985), to perform back-transformation from the
 172 Box-Cox scale. However, for convenience, the empirical exponential regression relationship
 173 ($y' = A + B R^x$) may be used to convert Box-Cox scale quantity x to the original scale
 174 quantity y' and subtracting any constant applied if necessary. This tends to be a good
 175 nonlinear relationship for CO₂ flux (see Figure 2) and for N₂O flux.

176

177 The median value of a sample can be a robust statistic but it will be influenced by the
 178 presence of a heavy- or long-tail. However, there are many heavy-tailed distributions such as
 179 extreme value Type I (Gumbel) and extreme value Type II (Weibull) among others. These

180 Type I and II along with Type III (Frèchet) are special cases of the GEV distribution (Coles,
 181 2001). Rather than focussing on these three special cases, GEV is used generally in the data
 182 analysis presented here. The shape parameter η allows fitting of this distribution to a variety
 183 of data histogram shapes. In heavy-tailed distributions the sample mean is pulled away from
 184 the majority of the data values and can be greatly overestimated. Fitting the GEV ensures that
 185 the mean estimate represents the majority of the data and thus mitigates overestimation bias.
 186 From the estimate of the GEV shape parameter η , one can see which thick-tail type
 187 distribution describes the data best. A positive shape parameter indicates the Frèchet
 188 distribution, zero the Gumbel, and negative the Weibull (Eastoe, 2017). A more important
 189 outcome is the estimate of data average (μ) that is relatively free of the effect of data in the
 190 long tail. The estimate of μ is calculated as a function of the GEV parameters η and α (eq. 3
 191 and eq. 4). Parameter α is the position of the majority of data peak similar to the geometric
 192 mean or mode position in skew distributions.

193

194 GEV is a simple three parameter distribution with cumulative distribution function, $F(x)$, and
 195 probability density function, $f(x)$, defined as follows.

196 The cumulative distribution function for the GEV (Smith, 1986; Martins et al. 2017) is given
 197 by

$$F(x) = \exp \left\{ - \left[1 + \eta \frac{(x - \xi)}{\alpha} \right]^{\frac{1}{\eta}} \right\} \text{ when } \eta \neq 0$$

$$198 \quad F(x) = \exp \left\{ - \exp \left[- \frac{(x - \xi)}{\alpha} \right] \right\} \text{ when } \eta = 0 \quad \text{eq. 3}$$

199 with ξ being the data peak location parameter, η the shape parameter and α the scale
 200 parameter.

201 The formula for probability density function (Singh, 1998) is

$$202 \quad f(x) = \frac{1}{\alpha} \left[1 - \frac{\eta}{\alpha} (x - \xi) \frac{1 - \eta}{\eta} \right] \exp \left[- \left(1 - \frac{\eta}{\alpha} (x - \xi) \right) \right]^{\frac{1}{\eta}} \quad \text{eq. 4}$$

203 From these references the mean (μ), standard deviation (σ) and skewness (γ) of the GEV
 204 distribution can be calculated:

$$\hat{\xi} = \hat{\mu} - \frac{\hat{\sigma}}{\hat{\eta}} \{1 - \Gamma(1 + \hat{\eta})\}$$

$$\hat{\sigma} = \frac{\hat{\sigma} \hat{\eta}}{\{\Gamma(1 + 2\hat{\eta}) - [\Gamma(1 + \hat{\eta})]^2\}^{1/2}}$$

$$205 \hat{\gamma} = \text{sign}(\hat{\eta}) \frac{-\Gamma(1+3\hat{\eta})+3\Gamma(1+\hat{\eta})\Gamma(1+2\hat{\eta})-2[\Gamma(1+\hat{\eta})]^3}{\{\Gamma(1+2\hat{\eta})-[\Gamma(1+\hat{\eta})]^2\}^{3/2}} \text{ eq. 5}$$

206 Symbol Γ denotes the gamma function.

207 The quantile function [the inverse of $F(x)$] of the GEV distribution is:

$$208 F^{-1}(u) = \zeta - \frac{\alpha}{\eta} [1 - \{-\log u\}^{-\eta}] \text{ with } 0 < u < 1 \quad \text{eq. 6}$$

209 When the interest is to estimate re-occurrence of (say) the maximum of a particular pollutant,
 210 then the value $F^{-1}(1 - u)$ is the return level associated with the return period $1/u$ (Sexto et
 211 al. 2013).

212

213 3. RESULTS AND DISCUSSION

214

215 The nature of GHG data is such that any spot value may not be representative of the flux size
 216 in a particular agricultural field. This is why the data in this study was collected every day
 217 over a period of 47 d. However, this extra time dimension creates a need to summarise data
 218 so the treatments may be compared by simple analysis of variance (ANOVA) based on the
 219 statistical design. Alternatively, a repeated measurement ANOVA of the design may be
 220 carried out without summarizing the data and a simple randomised block ANOVA using
 221 meaningful summary statistics is also desirable. For this purpose, the time course profile may
 222 be modelled if a suitable model is identifiable, otherwise calculating the area under the curve
 223 can be a good surrogate summary.

224

225 To understand the averaging problem, the sample average (implicitly assuming normal
226 distribution), log-normal based mean, Box-Cox transformation based mean and mean from
227 the fit of GEV distribution were considered. This exercise was completed with CO₂ flux data
228 (Tables 3 and 4). The N₂O flux data (Tables 5 and 6) had very small scale size observations
229 and both positive and negative values. Thus, the algorithm to estimate λ did not converge to a
230 satisfactory solution. Even to calculate the mean via the log-scale, it was necessary to use
231 $\ln(X + c)$ with $c = \text{absolute value of the minimum (N}_2\text{O flux)} + 1.0$. Because of this difficulty
232 the results for N₂O flux amended as above are included. From the parameters of the GEV
233 distribution, the GEV mean for CO₂ flux was estimated to be 60.14 kg CO₂-C ha⁻¹ d⁻¹ shown
234 in Table 7 (note that $(\mu - \xi)$ is the contribution from data in the heavy tail). Similarly, the
235 estimated GEV mean for N₂O flux (net of the added constant of 1.02029) was 0.01571 kg
236 N₂O-N ha⁻¹ day⁻¹ ($\{1.036 - 1.02029\}$; Table 8).

237

238 3.1. Carbon dioxide flux data

239 The example data employed here demonstrate a heavy tailed distribution as shown by
240 skewness of 1.688 ± 0.104 and kurtosis of 5.555 ± 0.208 due to the presence of excessive
241 hot-moments or extreme values (see Figure 3). This feature of data distributions means non-
242 robust statistics such as the arithmetic mean will be biased positively. When examined on the
243 log-scale the data were still non-normal. Similarly, data on the Box-Cox Scale with $\lambda = 0.1$
244 did not become normal. However, when using the generalised extreme value distribution, the
245 data were found to be consistent with that distribution (Figure 4).

246

247 When analysing CO₂ flux data, the generalised extreme value distribution fitted successfully
248 to individual treatments and overall (GEV parameters given in Table 4) and it provided a
249 better description of the data compared to the normal, log-normal and Box-Cox transformed

250 data. It therefore seems GEV is a viable option to analyse long-tailed or heavy-tailed GHG
251 data. The analysis of CO₂ data shows that fitting the GEV distribution can reduce bias from
252 the sample mean estimate (Table 7) and the standard deviation is also lower. From the fitted
253 parameters of the GEV distribution (Table 4), the mean μ and standard deviation σ were
254 calculated (Table 7).

255

256 *3.2. Nitrous oxide flux data*

257 Nitrous oxide flux data appear very different across the 12 plots in this study. Values range
258 from high positive values to negative values (Figure 5). Thus, the data for N₂O flux were
259 much more complex than CO₂ flux data due to the presence of negative fluxes. These data
260 form mixtures of distributions. Even GEV did not fit to the individual plot data sets entirely
261 satisfactorily (Figure 6) despite the addition of a constant of 1.02029 (i.e. 1 + minimum
262 absolute data value) to the data. From the GEV parameter estimates in Table 6, the estimates
263 of mean μ and standard deviation σ were calculated (Table 8).

264

265 **CONCLUSIONS**

266

267 This study shows that when analysing GHG data from agricultural fields, detailed analysis is
268 required before proceeding to the application of a suitable methodology. Black-box or default
269 statistics such as simple sample mean can give biased estimates. It is prudent to test implicit
270 distributional assumptions in order to identify an appropriate methodology.

271 From the above, we can draw up a general strategy for GHG field data analysis:

- 272 1. Check the distribution of the data and see if it is normally distributed.
- 273 2. Check for presence of hot-moments or outliers and deal with them if present (see Dhanoa
274 et al. 2016 for various tests).

275 3. If the data distribution appears to be skew, consider if the data are expected to follow a
276 log-normal distribution. Check the distribution after logarithmic transformation. If data
277 observations include negative and/or zero values, then $\ln(x + c)$ should be used with value of
278 constant c such that all data observations are positive. As explained above, when converting
279 any log-scale statistics on to the original scale the Finney (1941) correction must be applied
280 (Dhanoa, 2017) and any constant that was added must be subtracted.

281 4. If the majority of the data appear to be normal apart from a few outliers, then the Box-
282 Cox transformation may be considered. When $\lambda = 0.0$, logarithmic transformation is indicated
283 otherwise use the Box-Cox scale as described above. Again, add a constant c to make all data
284 positive.

285 5. However, if the distribution tail is long with many observations in it, then the option of
286 the generalised extreme value distribution may be relevant.

287

288 In the case of GHG data studies, many of the results may end up in national inventories.

289 Thus, it is necessary and important to follow all procedures that help minimise any biases in
290 the data summaries, to enable meta-analysis and other statistical comparisons of treatments
291 and studies provide suitable measure(s) of uncertainty.

292

293 **ACKNOWLEDGEMENTS**

294

295 The work was supported by the Biotechnology and Biological Sciences Research Council
296 (BB/P01268X/1, BBS/E/C/000I0320).

297

298 **REFERENCES**

299

- 300 Atkinson, A.C. (1982) Regression diagnostics, transformations and constructed variables
301 (with discussion). *Journal of the Royal Statistical Society, Series B*, 44, 1-36.
302
- 303 Atkinson, A.C. (1985) *Plots, Transformations and Regression*. Oxford University Press,
304 Oxford.
305
- 306 Battista, G., Pagliaroli, T., Mauri, L., Basilicata, C., and de Lieto Vollaro, R. (2016)
307 Assessment of the Air Pollution Level in the City of Rome (Italy). *Sustainability*, 8 (9): 838-
308 852
309
- 310 Beniston, M. (2004). The 2003 heat wave in Europe: A shape of things to come? An analysis
311 based on Swiss climatological data and model simulation. *Geophysical Research Letters*, 31,
312 L02202
313
- 314 Box, G.E.P. and Cox, D.R. (1964) An analysis of transformations (with discussion). *Journal*
315 *of the Royal Statistical Society, Series B*, 26, 211-46.
316
- 317 Coles, S. (2001) *An Introduction to Statistical Modeling of Extreme Values*. Springer-Verlag,
318 London, 224pp.
319
- 320 Dhanoa, M.S., Sanderson, R., Lopez, S., Kebreab, E. and France, J. (2016) Consequences of
321 metabolic scaling and log-scale allometry on means and variances and parameter estimates
322 from Type I and Type II linear regression models. *e-Planet* 14(1), 1–10.
323
- 324 Dhanoa, M.S. (2017) Cinderella’s transformation. *Significance*, 14(2), 46.

325

326 Dixon, E. R., Blackwell, M. S. A., Dhanoa, M. S., Berryman, Z., de la Fuente Martinez, N.,
327 Junquera, D., Martinez, A., Murray, P. J., Kemp, H. F., Meier-Augenstein, W., Duffy, A. &
328 Bol, R. (2010) Measurement at the field scale of soil $\delta^{13}\text{C}$ and $\delta^{15}\text{N}$ under improved
329 grassland. *Rapid Communications in Mass Spectrometry* 24, 511-518.

330 <http://dx.doi.org/10.1002/rcm.4345>

331

332 Eastoe, E. (2017) Extreme value distributions. *Significance*, 14, 12–13, doi:10.1111/j.1740-
333 9713.2017.01014.

334

335 Edwards W, Moles AT, Chong C. (2015) Generalised extreme value distributions provide a
336 natural hypothesis for the shape of seed mass distributions. *PLoS One*. 10(4), e0121724.

337 doi:10.1371/journal.pone.0121724

338

339 Ercelebi, S. G. and Toros, H. (2009) Extreme Value Analysis of Istanbul Air Pollution Data.
340 *CLEAN - Soil Air Water* 37(2):122–131

341

342 Evans, M., Hastings, N.A.J. and Peacock, J.B. (2000) *Statistical Distributions* (3rd Edition).

343 John Wiley & Sons, New York.

344

345 Finney, D.J. (1941) On the distribution of a variate whose logarithm is normally distributed.

346 *Journal of the Royal Statistical Society Series B*, 7, 155-161.

347

- 348 Fisher, R.A. and Tippett, L.H.C. (1928) Limiting forms of the frequency distributions of the
349 largest or smallest member of a sample. *Proceedings of the Cambridge Philosophical*
350 *Society* 24, 180–90.
- 351
- 352 Katz RW (2010) Statistics of extremes in climate change. *Climate Change*, 100, 71–76.
- 353
- 354 Kharin, V.V. and Zwiers, F.W. (2004) Estimating Extremes in Transient Climate Change
355 Simulations. *Journal of Climate*, 18, 1156-1173. <https://doi.org/10.1175/JCLI3320.1>.
- 356
- 357 Korkmaz, M. C. (2015) Two-sided generalised Gumbel distribution with application to air
358 pollution data. *International Journal of Statistical Distributions and Applications*, 1 (1), 19-
359 26.
- 360
- 361 Küchenhoff, H. and Thamerus, M. (1996) Extreme value analysis of Munich pollution data.
362 *Environmental and Ecological Statistics*, 3 (2), 127-141.
- 363
- 364 Loick, N., Dixon, E., Abalos, D., Vallejo, A., Matthews, P., McGeough, K., Watson, C.,
365 Baggs, E., Cardenas, L.M. (2017) “Hot spots” of N and C impact nitric oxide, nitrous oxide
366 and nitrogen gas emissions from a UK grassland soil. *Geoderma*, 305 (2017) 336–345
- 367
- 368 Louro, A., Sawamoto, T., Chadwick, D., Pezzolla, D., Bol, R., Baez, D. and Cardenas, L.
369 (2013) Effect of slurry and ammonium nitrate application on greenhouse gas fluxes of a
370 grassland soil under atypical South West England weather conditions. *Agriculture,*
371 *Ecosystems and Environment*, 181, 1-11.
- 372

- 373 Martins, L. D., Wikuats, C. F. H., Capucim, M. N., de Almeida, D. S., da Costa, S. C.,
374 Albuquerque, T., Carvalho, V. S. B., de Freitas, E. D., Andrade, M. de F., Martins, J. A.
375 (2017). Extreme value analysis of air pollution data and their comparison between two large
376 urban regions of South America. *Weather and Climate Extremes*, **18**, (December), Pages 44-
377 54
378
- 379 McFadden, Daniel (1978). Modeling the Choice of Residential Location. *Transportation*
380 *Research Record* 673, 72–77.
381
- 382 Pezzolla, D., Bol, R., Gigliotti, G., Sawamoto, T., Louro-López, A., Cardenas, L. and
383 Chadwick, D. (2012) Greenhouse gas (GHG) emissions from soils amended with digestate
384 derived from anaerobic treatment of food waste. *Rapid Communications in Mass*
385 *Spectrometry*, 26, 2422-2430.
386
- 387 Sexto, B. M., Vaquera, H.H. and Arnold, B. C. (2013) Use of the Dagum distribution for
388 modelling tropospheric Ozone levels, *Journal of Environmental Statistics*, 5 (5), 1-11.
389
- 390 Sheather, S.J. and Jones, M.C. (1991) A reliable data-based bandwidth selection method for
391 kernel density estimation. *Journal of the Royal Statistical Society, Series B*, 53, 683-690.
392
- 393 Singh, V.P. (1998) Generalized extreme value distribution. In *Entropy-Based Parameter*
394 *Estimation in Hydrology*, vol. 30, pp169-183. Springer, Dordrecht, the Netherlands.
395

Figure 1. Kernel density plot of (a.) CO₂ and (b.) N₂O flux data showing main area of data frequency and the composition of the heavy-tail.

Figure 2. Empirical relationship ($y' = a + B R^x$) between CO₂ flux data on the original scale and the corresponding values scaled according to the Box-Cox (1964) transformation with $\lambda = 0.1$

Figure 3. CO₂ flux data showing observations contributing to skewness and heavy tail.

Figure 4. Probability plot when modelling CO₂ flux data using the GEV distribution.

Figure 5. N₂O flux data showing observations that contribute to skewness and a heavy tail.

Figure 6. Probability plot when modelling N₂O flux data using the GEV distribution.

1 Table 1. Chemical composition of applied slurry and digestate

2

Property	Units	Slurry application	Digestate application
Dry matter	%	6.5	4.8
Density	kg l ⁻¹	1.006	1.00
Ammonium, NH ₄ ⁺ N	g kg ⁻¹ dry matter	18.5	97.3
Nitrate, NO ₃ ⁻ N	g kg ⁻¹ dry matter	0.0	0.0
Total N	% of dry matter	2.67	16.9
pH	-	7.3	8.16
Total carbon	% of dry matter	38.4	38.6
C:N ratio	-	14.4	2.3

3

4

5

6

7

8

9

10

11

12

13

14

15

16

17

18

19

20

21

22

23

24
25
26

Table 2. Glossary of input parameters calculated from GHG flux data

GHG data Parameter	Description	References
AM	arithmetic sample mean	
γ	skewness	
$\hat{\mu}$	sample mean on the log-scale	Finney (1941)
$\hat{\sigma}^2$	sample variance on the log-scale	Finney (1941)
y	Original data value	Box and Cox (1964)
c	Positive constant which when added to all dataset values makes all data above zero (only if negative values exist)	Box and Cox (1964)
λ	transformation parameter to fit normal distribution curve	Box and Cox (1964)
η	data peak location parameter	GEV, Smith, 1986
ξ	data distribution shape parameter	GEV, Smith, 1986
α	scale parameter	GEV, Smith, 1986
Γ	Gamma function	Singh, 1998
u	Return period of peak level	Sexto et al. 2013

27
28
29
30
31
32
33
34
35
36
37
38
39

40

41 Table 3. Sample statistics for CO₂ flux data (kg CO₂-C ha⁻¹ d⁻¹) over 47 d from three replicate
 42 plots of each of four treatments comprising digestate (DG), cattle slurry (SL), control (CN)
 43 and ammonium nitrate (AN).

44

Treatment-Plot	Sample Mean	Sample SD	Skewness	Kurtosis
DG-1	82.71	58.860	2.071	4.383
DG-2	85.67	44.654	0.400	-0.875
DG-3	65.52	30.580	0.711	-0.506
SL-1	50.56	15.525	0.465	-0.597
SL-2	70.34	36.105	0.882	-0.394
SL-3	73.01	28.340	0.773	-0.171
CN-1	56.85	33.904	1.532	2.377
CN-2	58.19	22.458	1.094	1.007
CN-3	35.07	18.287	1.431	2.540
AN-1	70.24	31.530	0.375	-0.645
AN-2	67.24	18.762	-0.139	-0.379
AN-3	72.07	33.038	0.687	0.406
Overall	Mean = 65.62 Median = 58.39	35.399	1.688	5.555

45

46

47 Table 4. Mean plot values for CO₂ flux data (kg CO₂-C ha⁻¹ d⁻¹) over 47 d from three
 48 replicate plots of each of four treatments comprising digestate (DG), cattle slurry (SL),
 49 control (CN) and ammonium nitrate (AN) estimated on the log-scale, Box-Cox scale (Box
 50 and Cox, 1964) and by generalised extreme value analysis (GEV).

51

Treatment-Plot	Mean log- scale*	Mean Box- Cox scale**	GEV		
			ξ	η	α
DG-1	81.58	70.49	54.20	0.371	26.332
DG-2	87.21	74.93	65.55	-0.050	37.020
DG-3	65.68	59.71	48.60	0.260	19.977
SL-1	50.64	48.51	43.78	-0.075	13.061
SL-2	70.38	63.12	50.66	0.291	21.909
SL-3	73.13	68.51	59.54	0.054	21.171
CN-1	56.70	49.93	39.89	0.265	18.901
CN-2	58.27	54.76	48.15	0.016	16.890
CN-3	35.86	31.19	27.19	0.011	13.504
AN-1	71.03	63.75	57.04	-0.123	27.561
AN-2	67.67	64.60	61.17	-0.338	19.038
AN-3	72.82	65.44	57.46	-0.048	27.178
Overall	65.89	58.53	48.94	0.116	23.670

52

53 * Finney (1941) correction applied.

54 ** Transformed back from Box-Cox scale using an empirical regression relationship,
 55 viz. $y' = a + B R^x$ where $y = \text{CO}_2$ flux and $x = \text{flux on Box-Cox scale with } \lambda = 0.1$.

56

57

58 Table 5. Sample statistics for N₂O flux data (kg N₂O-N ha⁻¹ d⁻¹) over 47 d from three
 59 replicate plots of each of four treatments comprising digestate (DG), cattle slurry (SL),
 60 control (CN) and ammonium nitrate (AN).

61

Treatment-Plot	Sample Mean	Sample SD	Skewness	Kurtosis
DG-1	1.045	0.0277	2.391	6.013
DG-2	1.046	0.0185	1.206	1.464
DG-3	1.036	0.0120	1.020	3.140
SL-1	1.031	0.0047	-0.152	-0.524
SL-2	1.051	0.0391	1.982	3.023
SL-3	1.038	0.0123	-0.312	1.255
CN-1	1.032	0.0108	2.007	4.709
CN-2	1.033	0.0082	1.062	1.984
CN-3	1.026	0.0050	-1.435	3.065
AN-1	1.042	0.0165	1.137	0.326
AN-2	1.033	0.0081	0.181	-0.473
AN-3	1.039	0.0126	1.562	2.538
Overall	Mean = 1.038 Median = 1.033	0.01868	3.457	1.766

62

63

64 Table 6. Mean plot values for N₂O flux data (kg N₂O-N ha⁻¹ d⁻¹) over 47 d from three
 65 replicate plots of each of four treatments comprising digestate (DG), cattle slurry (SL),
 66 control (CN) and ammonium nitrate (AN) estimated on the log-scale, Box-Cox scale (Box
 67 and Cox, 1964) and by generalised extreme value analysis.

68

Treatment-Plot	Mean log- scale*	Mean Box- Cox scale**	GEV		
			ξ	η	α
DG-1	1.0450	1.059	1.033	0.278	0.0124
DG-2	1.0458	1.061	NA	NA	NA
DG-3	1.0358	1.053	1.031	-0.112	0.0107
SL-1	1.0307	1.049	NA	NA	NA
SL-2	1.0513	1.064	NA	NA	NA
SL-3	1.03767	1.055	1.034	-0.331	0.0128
CN-1	1.0320	1.050	1.027	0.089	0.0068
CN-2	1.0326	1.050	1.029	-0.062	0.0038
CN-3	1.0259	1.044	NA	NA	NA
AN-1	1.0421	1.058	NA	NA	NA
AN-2	1.0333	1.051	1.030	-0.242	0.0078
AN-3	1.0394	1.056	1.034	0.058	0.0086
Overall	1.0376	1.037	1.030	0.072	0.0108

69

70 * Finney (1941) correction applied.

71 ** Transformed back from Box-Cox scale using an empirical regression relationship, viz.

72 $y = A + B R^x$ where $y = \text{N}_2\text{O flux (with added constant)}$ and $x = \text{flux on Box-Cox scale with}$ 73 $\lambda = -4.0$.

74 NA = Not available, distribution did not fit.

75

76 Table 7. CO₂ flux sample mean and mean and standard deviation as calculated from the
 77 parameters of the GEV distribution.

78

Treatment-Plot	Sample Mean	GEV	
		$\hat{\mu}$	$\hat{\sigma}$
DG-1	82.71	62.06	25.245
DG-2	85.67	88.84	50.949
DG-3	65.52	55.94	20.197
SL-1	50.56	52.35	18.672
SL-2	70.34	58.26	21.773
SL-3	73.01	70.68	25.432
CN-1	56.85	46.78	19.057
CN-2	58.19	57.63	21.213
CN-3	35.07	34.84	17.081
AN-1	70.24	76.73	42.851
AN-2	67.24	81.60	53.569
AN-3	72.07	74.49	37.260
Overall	65.62	60.14	26.679

79

80

81 Table 8. N₂O flux sample mean and mean and standard deviation calculated from the
 82 parameters of the fitted GEV distribution (flux data used include the added constant 1.02029
 83 to overcome negative and zero values in the original data).

84

Treatment-Plot	Sample Mean	GEV	
		$\hat{\mu}$	$\hat{\sigma}$
DG-1	1.045	1.037	0.0124
DG-2	1.046	NA	NA
DG-3	1.036	1.038	0.0163
SL-1	1.031	NA	NA
SL-2	1.051	NA	NA
SL-3	1.038	1.047	0.0350
CN-1	1.032	1.031	0.0079
CN-2	1.033	1.034	0.0095
CN-3	1.026	NA	NA
AN-1	1.042	NA	NA
AN-2	1.033	1.037	0.0158
AN-3	1.039	1.038	0.0103
Overall	1.038	1.036	0.0127
		Net 0.01571	

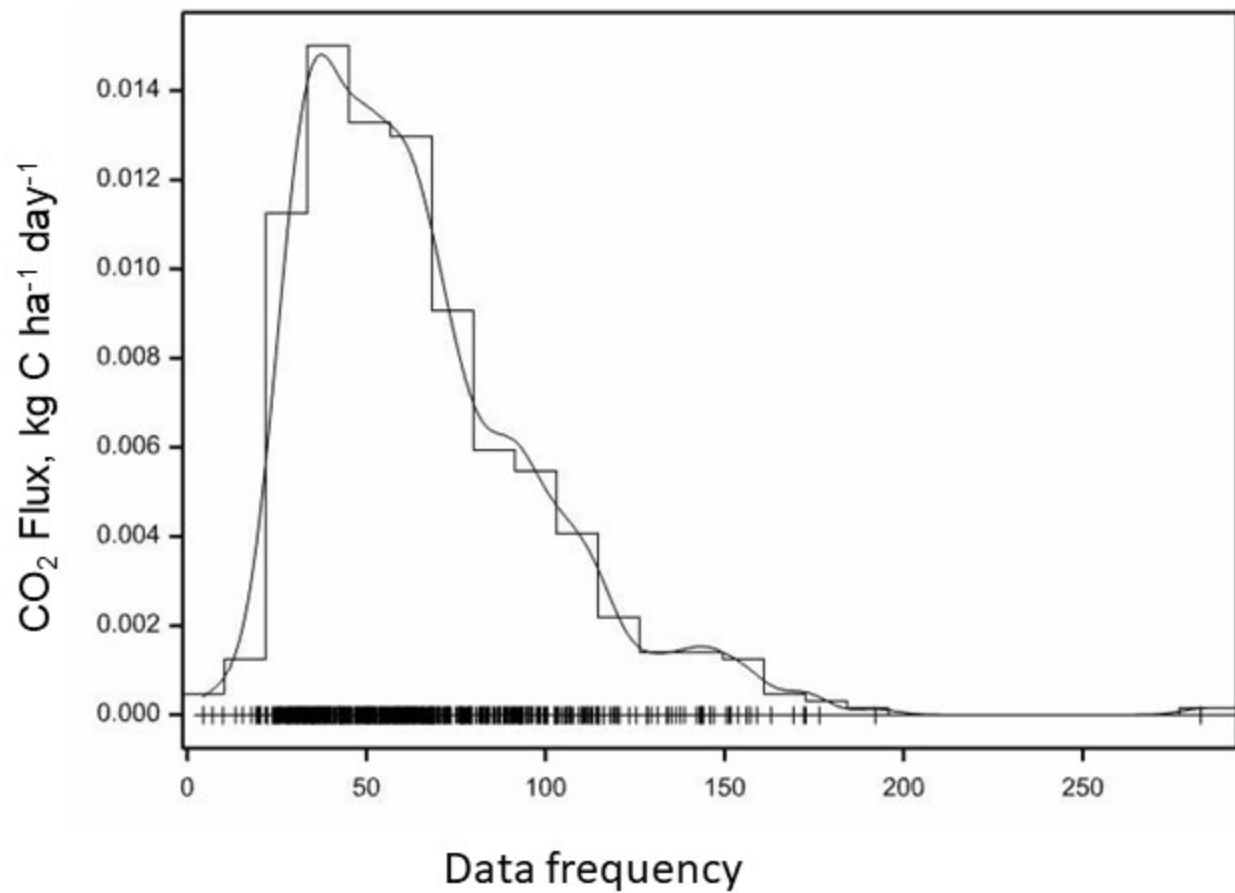
85

86 NA = Not available, distribution did not fit.

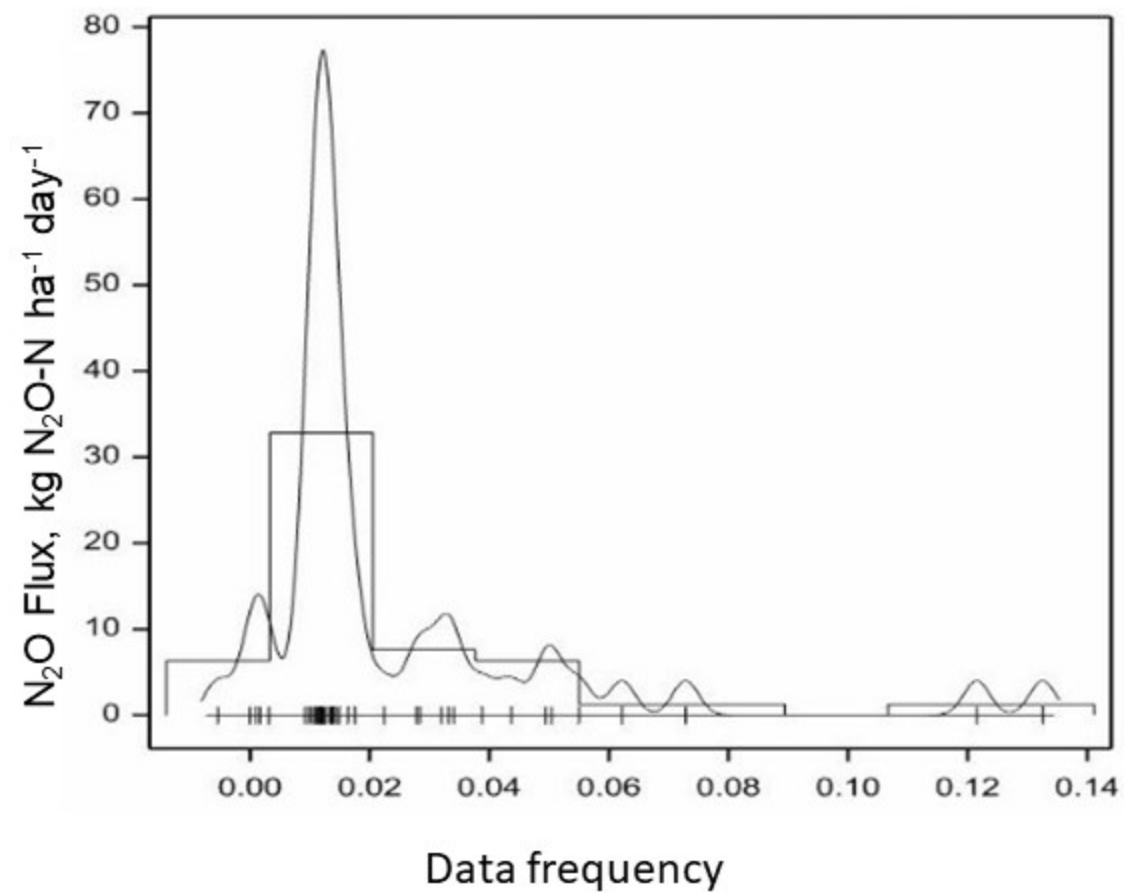
87

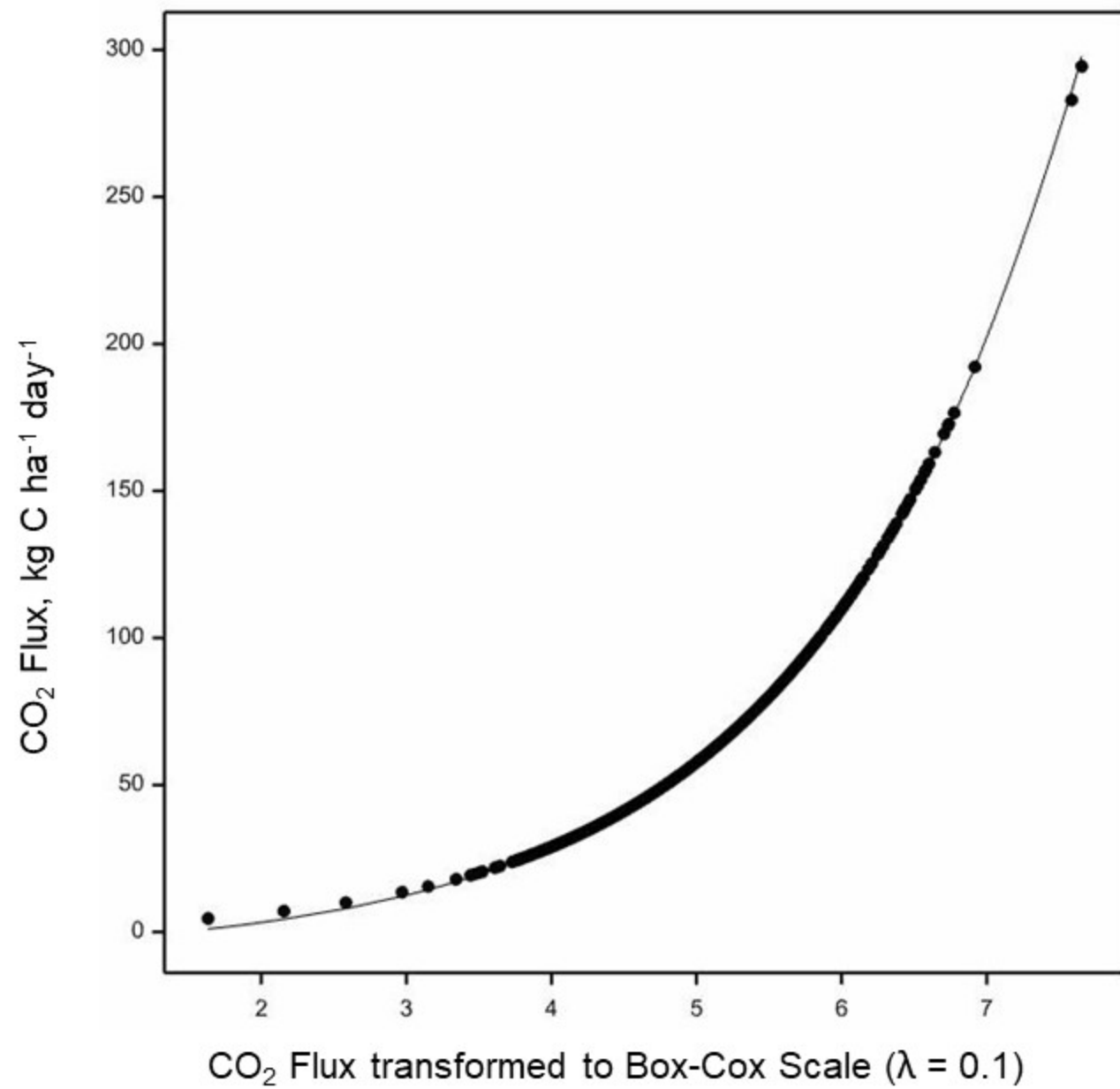
88

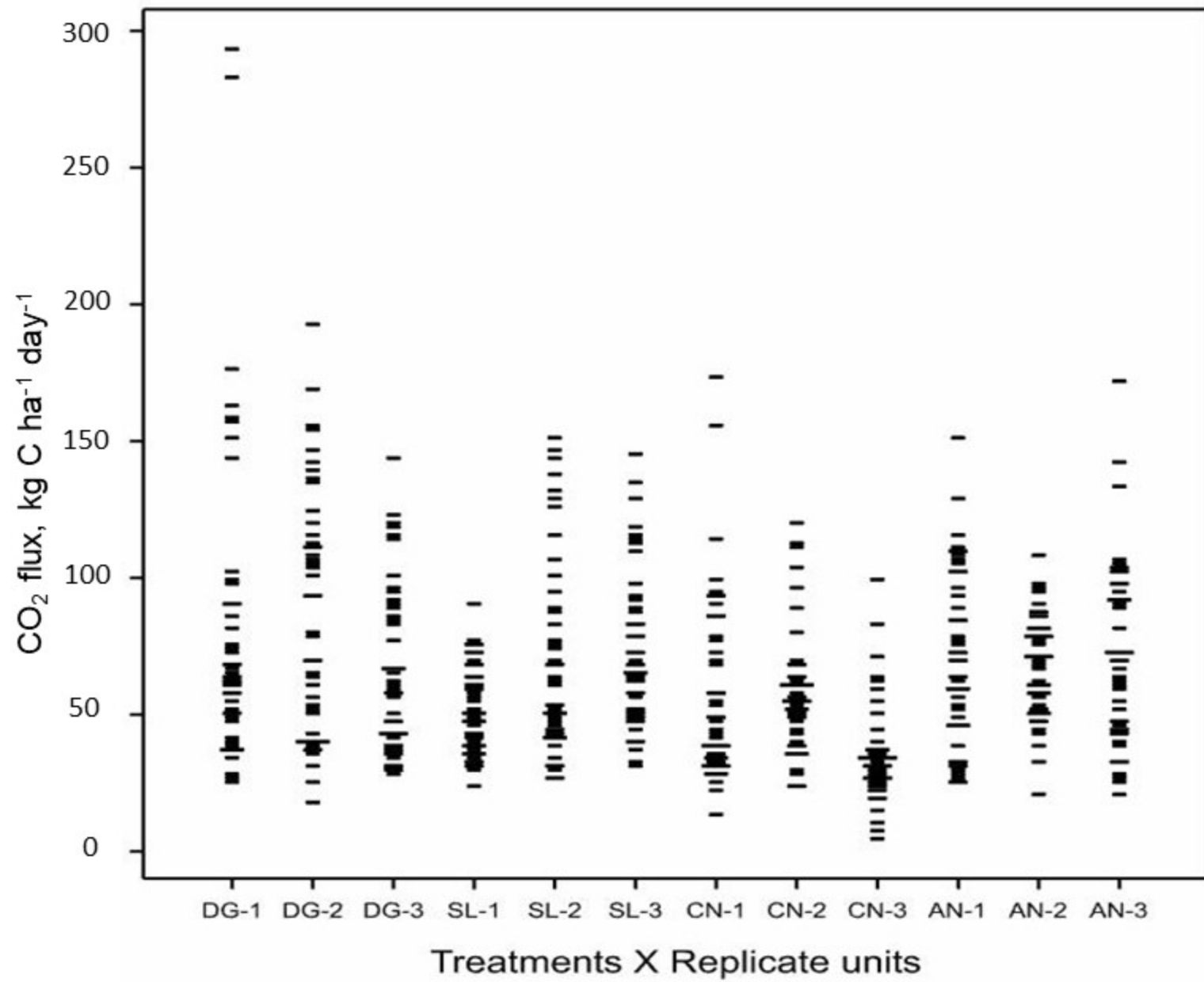
a.

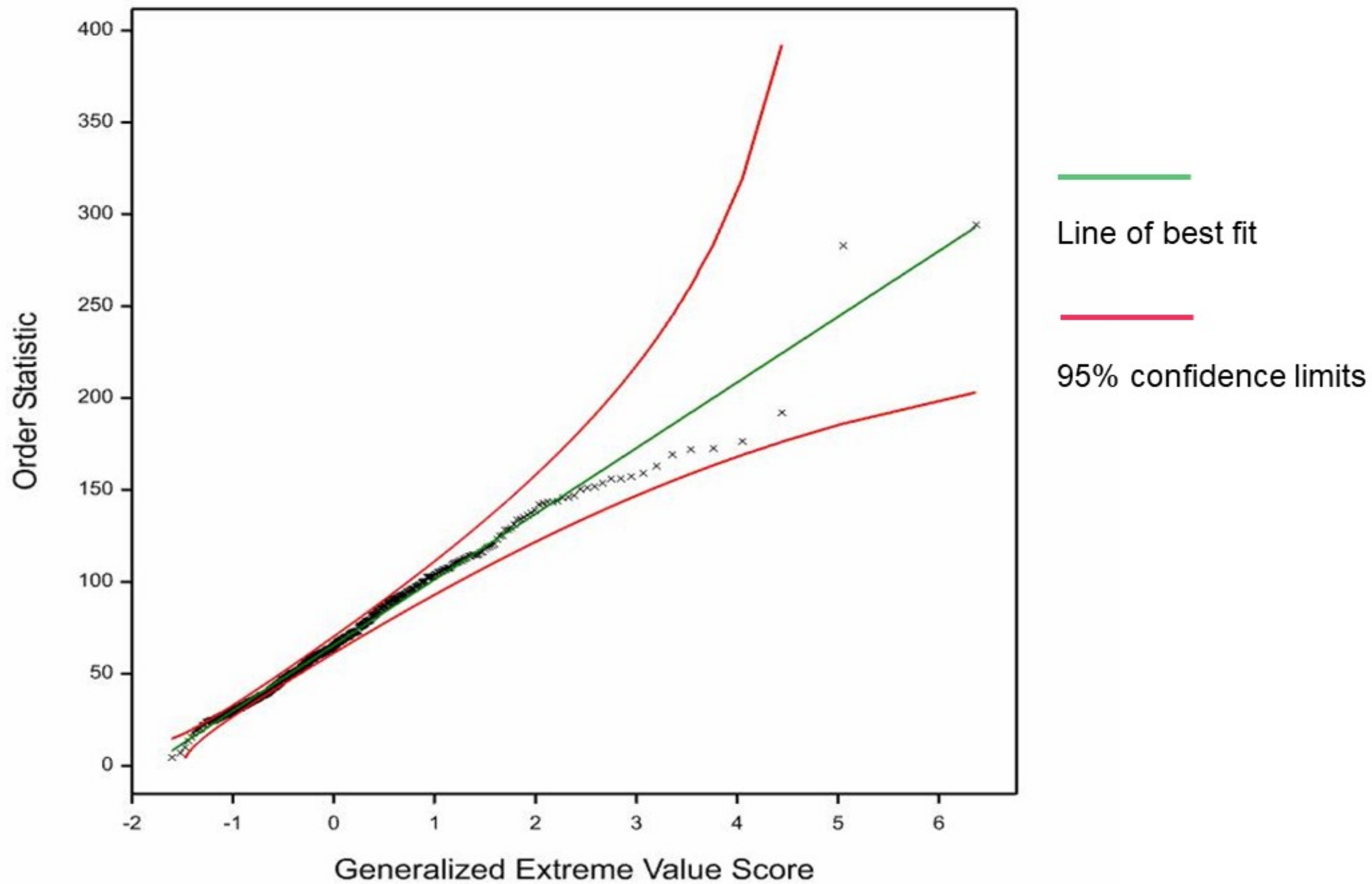


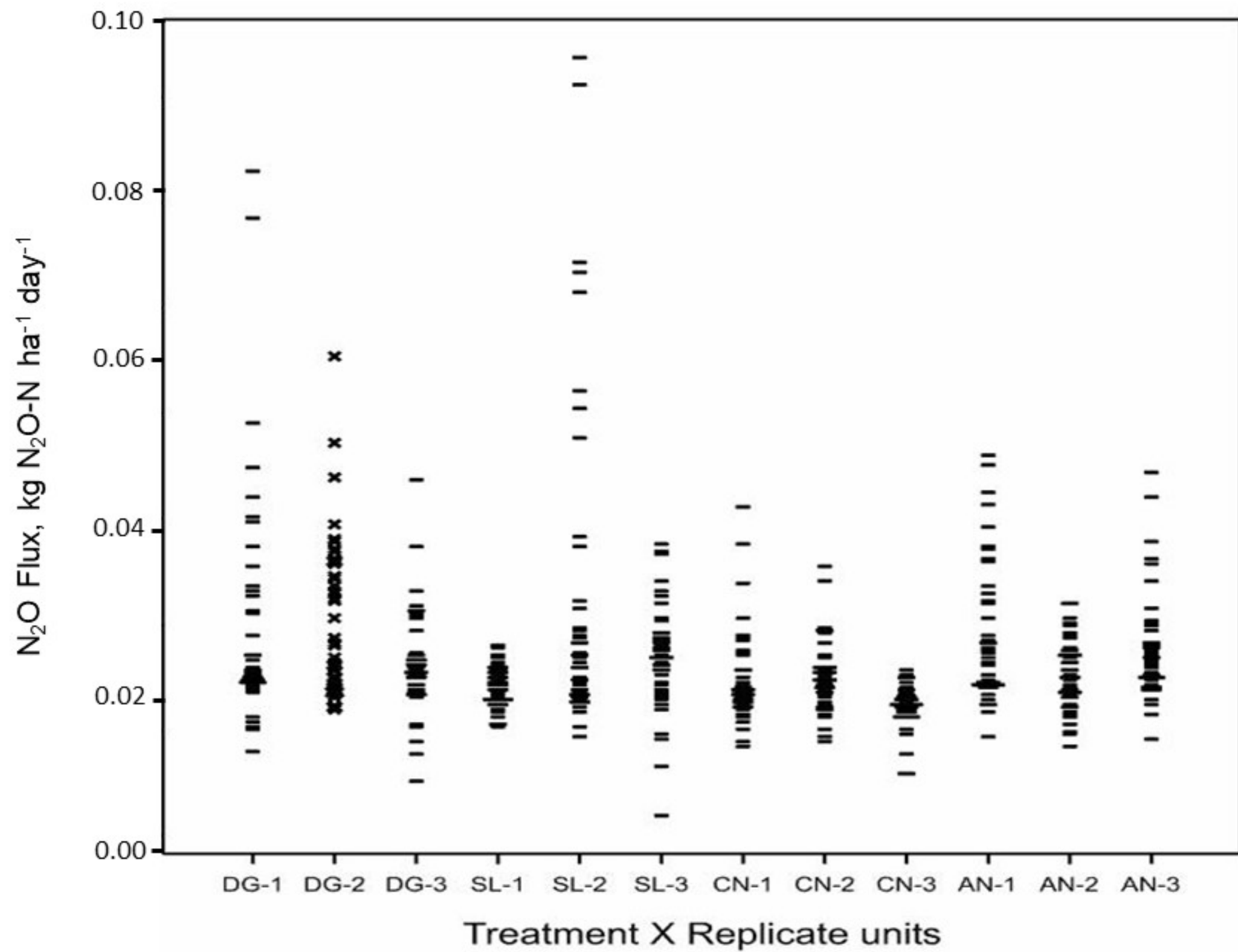
b.

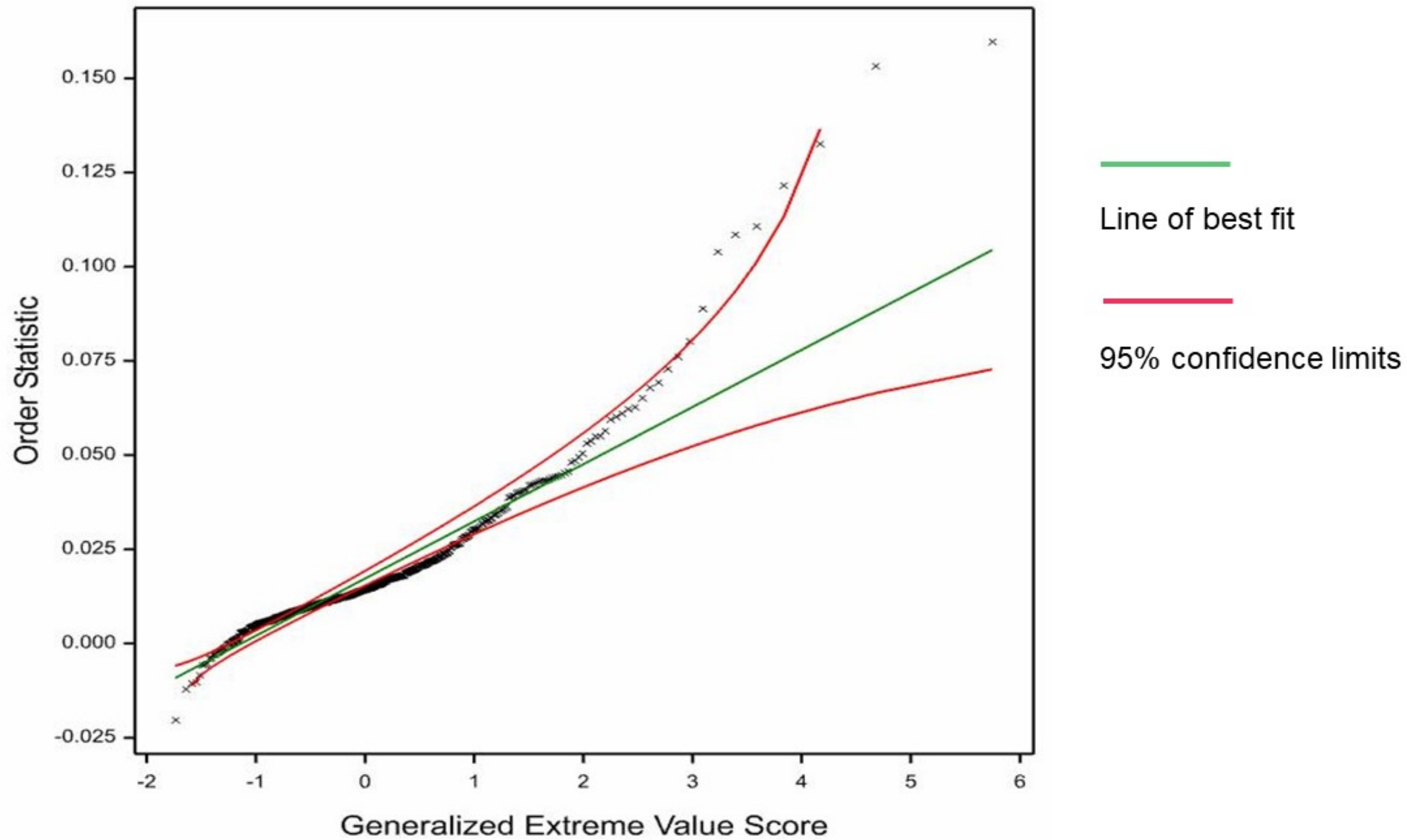












Highlights:

- Using sample means may be overestimating GHG fluxes
- GEV solves excessive skewness and kurtosis of greenhouse gas flux data
- Strategy of options for analysing GHG data rather than black-box approach
- CO₂ estimates from GEV less affected by data in the long tail than sample mean
- CO₂ estimates from Box-Cox are more affected by long-tail data than from GEV

Declaration of interests

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

The authors declare the following financial interests/personal relationships which may be considered as potential competing interests: