

1 A NOTE ON THE ANALYTIC APPROXIMATION OF EXCEEDANCE
2 PROBABILITIES IN HETEROGENEOUS POPULATIONS

3 By

4 H. S. Battey

5 *Department of Mathematics, Imperial College London, SW7 2AZ, UK*

6 SUMMARY

7 This note derives analytic approximations to exceedance probabilities for order statistics
8 from two heterogeneous populations. A limitation of this approach is that it entails a special
9 condition that needs to be checked or justified on a case by case basis.

10 *Some key words:* Analytic approximations; Competing risks, Extremal-Types Theorem,
11 Probabilistic screening.

12 **1 Introduction**

13 Suppose that X_1, \dots, X_{N_X} are N_X independent copies of the random variable X , whose
14 distribution function is F_X , and Z_1, \dots, Z_{N_Z} are N_Z copies of Z , whose distribution function
15 is F_Z . Both F_X and F_Z are assumed continuous. It is not known which of the $N_X + N_Z$
16 realizations of these random variables are from the X and Z populations. The following
17 quantities arise in diverse contexts:

- 18 (i) the probability that the set of largest $s \geq N_X$ observations contains all N_X observa-
19 tions from the X population;
- 20 (ii) the probability that the minimum of the X population exceeds the r th largest obser-
21 vation from the Z population;
- 22 (iii) the probability that the maximum of the $N_X + N_Z$ observations belongs to the X
23 population.

24 Exact solutions for such probabilities are rarely available, so that reliance is often on nu-
25 merical approximations, making dependence on key aspects inexplicit.

26 The first two questions emerge naturally in screening-type problems. Suppose the X
 27 population represents individuals with a particular disease, which can be assessed accurately
 28 by a costly, or otherwise inconvenient, procedure. To better target this resource, an initial
 29 screening is performed. Thus in this context X and Z may represent, for instance, systolic
 30 blood pressure, concentrations of toxins or solutes in the blood, etc. for individuals in the
 31 diseased and healthy group. The first question above indicates the likely success of the
 32 screening at detecting all individuals in the diseased group. The second quantifies the
 33 number of false positives entailed in the identification of all N_X cases. Similar issues arise
 34 in statistical contexts, where it is often desirable to screen a large number of potentially
 35 explanatory variables, sometimes with a view to assessing causality more rigorously through
 36 a full factorial or other designed experiment.

37 A key unifying observation for addressing these questions is that there is no loss of
 38 generality by treating one of the two distributions as standard uniform. This is because the
 39 transformed random variable $U_i \triangleq 1 - F_Z(Z_i)$ is uniformly distributed on $[0, 1]$. Here and
 40 henceforth \triangleq means equality by definition. On defining $V_i \triangleq 1 - F_Z(X_i)$ and modelling the
 41 density function of this random variable by $(1 - \gamma)v^{-\gamma}$ for $\gamma < 1$, $0 \leq v \leq 1$, an approach
 42 that would require careful justification in any particular context, probability calculations
 43 of the nature outlined above can be expressed in terms of beta integrals, approximable in
 44 terms of elementary functions using Stirling's formula.

45 In the expression $(1 - \gamma)v^{-\gamma}$, $\gamma = 0$ recovers the uniform density function and $\gamma \rightarrow 1$
 46 and $\gamma \rightarrow -\infty$ represent strong departure from the uniform distribution in both directions.

47 Of course, $1 - F_Z(Z_i)$ is not the only transformation that delivers uniform random
 48 variables and there are settings in which one of the other three possibilities, $F_Z(Z_i)$, $1 -$
 49 $F_X(X_i)$, $F_X(X_i)$, may be more fruitful. We discuss this choice in the context of the examples
 50 given.

51 A referee has pointed out work by Bairamov and Parsi (2011) and Bayramoglu and
 52 Eryilmaz (2015), who study a very similar problem. Their results apply under weaker con-
 53 ditions than in the present paper, at the expense of more complicated analytic expressions.
 54 In particular, X_1, \dots, X_{N_X} and Z_1, \dots, Z_{N_Z} are treated as exchangeable random variables,
 55 independent of one another in the first paper and dependent in the second. While these

56 two papers are the most closely related to the present work, there is an extensive literature
57 studying the exact and asymptotic distributions of the number of exceedances based on
58 order statistics. Notable early examples are Gumbel and von Schelling (1950) and Sarkadi
59 (1957) who consider exceedance probabilities of order statistics from two samples of poten-
60 tially different sizes drawn from the same population (i.e. $F_X = F_Z$). Certain special cases
61 are recoverable both from their calculations and ours, as indicated below.

62 **2 Exceedance probability formulae: examples**

63 **2.1 Competing maxima**

64 We first address the most challenging of the questions specified in §1. This serves as an
65 exemplar for the other cases. The probability that the maximum of $N_X + N_Z$ observations
66 belongs to the X population is the probability that $X_{\max} \triangleq \max\{X_1, \dots, X_{N_X}\}$ exceeds
67 $Z_{\max} \triangleq \max\{Z_1, \dots, Z_{N_Z}\}$. Since distribution functions are monotonically increasing,

$$p = \text{pr}(X_{\max} > Z_{\max}) = \text{pr}(V_{\min} < U_{\min}),$$

68 where $V_{\min} \triangleq \min\{V_1, \dots, V_{N_X}\}$ and $U_{\min} \triangleq \min\{U_1, \dots, U_{N_Z}\}$. Note that the density
69 function of U_{\min} at u is $N_Z(1-u)^{N_Z-1}$. Thus consider initially the event $A_i \triangleq \{V_i < U_{\min}\}$,
70 with associated probability

$$\text{pr}(A_i) = N_Z \int_0^1 v^{1-\gamma}(1-v)^{N_Z-1} dv.$$

71 This is of the form of a beta integral of indices $2 - \gamma$ and N_Z . Using $x\Gamma(x) = \Gamma(x+1)$ and
72 Stirling's formula in the form $\Gamma(x)/\Gamma(x+a) \simeq x^{-a}$ for large x and fixed a , we obtain

$$\text{pr}(V_i < U_{\min}) \simeq \frac{\Gamma(2-\gamma)\Gamma(N_Z+1)}{\Gamma\{N_Z+1+(1-\gamma)\}} \simeq \frac{\Gamma(2-\gamma)}{(N_Z+1)^{(1-\gamma)}}, \quad (N_Z \rightarrow \infty).$$

73 Since V_1, \dots, V_{N_X} are independent and identically distributed, for any fixed v ,

$$\text{pr}(V_{i_1} < v, \dots, V_{i_k} < v) = v^{k(1-\gamma)}$$

74 for an arbitrary set of k indices i_1, \dots, i_k . It follows that there are similar approximations
 75 to the probabilities of all combinations of joint events, specifically

$$\text{pr}(A_{i_1}, \dots, A_{i_k}) \simeq \frac{\Gamma\{k(1-\gamma) + 1\}}{(N_Z + 1)^{k(1-\gamma)}}, \quad (N_Z \rightarrow \infty).$$

76 We conclude that

$$p = \text{pr}\left(\bigcup_{i=1}^{N_X} A_i\right) \simeq \sum_{k=1}^{N_X} (-1)^{k-1} \binom{N_X}{k} \frac{\Gamma\{k(1-\gamma) + 1\}}{(N_Z + 1)^{k(1-\gamma)}}, \quad (N_Z \rightarrow \infty). \quad (1)$$

77 An analysis of convergence is given in the supplementary material. The approximation is
 78 essentially exact, the only error coming from the use of Stirling's formula. Indeed, Stirling's
 79 approximation to $\Gamma(x)$, while derived under the notional limiting operation $x \rightarrow \infty$ provides
 80 a remarkably accurate approximation even for small values of x . It is therefore reasonable
 81 to apply Stirling's formula to the numerators, giving the following approximation in terms
 82 of elementary functions:

$$p \approx \frac{\sqrt{2\pi}}{e} \sum_{k=1}^{N_X} (-1)^{k-1} \binom{N_X}{k} \frac{\{k(1-\gamma) + 1\}^{k(1-\gamma) + \frac{1}{2}}}{e^{k(1-\gamma)} (N_Z + 1)^{k(1-\gamma)}}. \quad (2)$$

83 Note that while the random variables V_1, \dots, V_{N_X} are independent, $\mathbb{I}\{V_1 < U_{\min}\}, \dots, \mathbb{I}\{V_{N_X} <$
 84 $U_{\min}\}$ are not, due to their mutual dependence on U_{\min} . It is for this reason that the full
 85 inclusion-exclusion formula is needed in (1) rather than the simplification

$$\text{pr}\left(\bigcup_{i=1}^n E_i\right) = \sum_{i=1}^n (-1)^{i-1} \binom{n}{i} \eta^i = 1 - (1 - \eta)^n,$$

86 applying to independent events E_i with probabilities $\eta = \text{pr}(E_i)$ for all i .

87 Reversing the roles of X_i and Z_i in the construction of the uniform random variables
 88 leads to a simpler formula. Let $\tilde{U}_i \triangleq F_X(X_i)$ and $W_i \triangleq F_X(Z_i)$ so that $\tilde{U}_1, \dots, \tilde{U}_{N_X}$ are
 89 independent uniformly distributed random variables. On modelling the density function
 90 of the random variables W_1, \dots, W_{N_Z} as $(1-\gamma)w^{-\gamma}$, the probability p is now $\text{pr}(\tilde{U}_{\max} >$
 91 $W_{\max})$, where $\tilde{U}_{\max} \triangleq \max\{\tilde{U}_1, \dots, \tilde{U}_{N_X}\}$ with density function $N_X u^{N_X-1}$ at u , and $W_{\max} \triangleq$

92 $\max\{W_1, \dots, W_{N_Z}\}$. The required probability is

$$p = \text{pr}\left(\bigcap_{i=1}^n \{W_i < \tilde{U}_{\max}\}\right) = N_X \int_0^1 w^{N_Z(1-\gamma)+N_X-1} dw = \frac{N_X}{N_X + N_Z(1-\gamma)}. \quad (3)$$

93 The choice between which transformation to use, leading to either (1) or (3) depends on
 94 which of F_Z and F_X is known, or for which of the populations V_1, \dots, V_{N_X} or W_1, \dots, W_{N_Z}
 95 the density parameterization $(1-\gamma)v^{-\gamma}$ is most reasonable.

96 If γ is set to zero in equation (1) and the binomial coefficient is written in terms of gamma
 97 functions and simplified using Stirling's approximation, the right hand side of equation (1)
 98 is, for $N_Z > N_X$,

$$\sum_{k=1}^{N_X} (-1)^{k-1} \left(\frac{N_X+1}{N_Z+1}\right)^k = \frac{N_X+1}{(N_X+1) + (N_Z+1)} \left\{ 1 + (-1)^{N_X+1} \left(\frac{N_X+1}{N_Z+1}\right)^{N_X+1} \right\}.$$

99 For large N_Z this is approximately $N_X/(N_X + N_Z)$ which is what one would obtain from
 100 direct calculation, noting that V_i and U_i are both uniformly distributed when $\gamma = 0$. This
 101 is equation (1.6) of Gumbel and von Schilling (1950). The formula (3) similarly becomes
 102 $N_X/(N_X + N_Z)$ when $\gamma = 0$. Also intuitively, for $\gamma = 1$, formula (1) reduces to

$$\sum_{k=1}^{N_X} (-1)^{k-1} \binom{N_Z}{k} = - \left\{ \sum_{k=0}^{N_X} \binom{N_Z}{k} - 1 \right\} = 1.$$

103 while for $\gamma \rightarrow -\infty$, it tends to zero. Similarly for formula (3).

104 2.2 Probabilities quantifying screening properties

105 As noted in §1, probabilities (i) and (ii) are relevant for assessing the properties of physical
 106 or abstract screening procedures. These questions lead to simpler calculations than that of
 107 §2.1 once the transformation to $U_i = 1 - F_Z(Z_i)$ and $V_i = 1 - F_Z(X_i)$ has been made.

108 We start by considering the probability that all members of the X population are dis-
 109 covered before the first falsely detected individual from the Z population. This is

$$\begin{aligned} \text{pr}(U_{\min} > V_{\max}) &= \text{pr}\left(\bigcap_{i=1}^{N_X} \{V_i < U_{\min}\}\right) \\ &= N_Z \int_0^1 v^{N_X(1-\gamma)} (1-v)^{N_Z-1} dv \simeq \frac{\Gamma\{N_X(1-\gamma) + 1\}}{(N_Z + 1)^{N_X(1-\gamma)}}, \quad (N_X \rightarrow \infty). \end{aligned}$$

110 The r th smallest order statistic $U_{(r)}$ constructed from U_1, \dots, U_{N_Z} is distributed as a
 111 beta random variable with indices r and $N_Z - r + 1$ so that its density function is given by

$$\frac{N_Z!}{(r-1)!(N_Z-r)!} v^{r-1} (1-v)^{N_Z-r}.$$

112 Probability (ii) is therefore addressed by

$$\begin{aligned} \text{pr}(U_{(r)} > V_{\max}) &= \frac{N_Z!}{(r-1)!(N_Z-r)!} \int_0^1 v^{N_X(1-\gamma)+r-1} (1-v)^{N_Z-r} dv \\ &= \frac{N_Z!}{(r-1)!} \frac{\Gamma\{N_X(1-\gamma) + r\}}{\Gamma\{N_X(1-\gamma) + N_Z + 1\}}, \end{aligned}$$

113 which can be further simplified using Stirling's formula. There are no technical difficulties
 114 in considering the generalization $\text{pr}(U_{(r)} > V_{(s)})$ for $1 \leq r \leq N_Z, 1 \leq s \leq N_X$.

115 3 Comparison to approximations based on limit theorems

116 In highly influential work, Fisher and Tippett (1928) characterized the set of probability
 117 laws \mathcal{L} such that, for independently distributed random variables X_1, \dots, X_n with a common
 118 but arbitrary distribution function F , there exist sequences $a_X(n)$ and $b_X(n)$ such that

$$\lim_{n \rightarrow \infty} \text{pr} \left(\frac{\max\{X_1, \dots, X_n\} - b_X(n)}{a_X(n)} \right) = \lim_{n \rightarrow \infty} F^n\{a_X(n)x + b_X(n)\} = L(x), \quad L \in \mathcal{L}. \quad (4)$$

119 When equation (4) holds, the random variable X is said to be in the domain of attraction
 120 of L . The set \mathcal{L} has just three elements. These are, in the notation of Fisher and Tippett
 121 (1928), for $\alpha > 0$:

$$\begin{aligned} \text{Type I} &: L(x) = \exp(-e^{-x}), \quad x \in \mathbb{R}; \\ \text{Type II} &: L(x) = L_\alpha(x) = \begin{cases} 0, & x \leq 0, \\ \exp(-x^{-\alpha}), & x > 0; \end{cases} \\ \text{Type III} &: L(x) = L_\alpha(x) = \begin{cases} \exp\{-(-x)^\alpha\}, & x < 0 \\ 1, & x \geq 0. \end{cases} \end{aligned}$$

122 Gnedenko (1943) proved that these are the only three types that can arise as limit laws. The
 123 scaling sequence $a_X(n)$ can be taken as 1 when X_1, \dots, X_n are in the domain of attraction
 124 of the Type I limit law.

125 These limit laws are plausible approximations for the context of §2.1 provided that
 126 both N_X and N_Z are large. In the context of §2.1, suppose that both types of random
 127 variables X_1, \dots, X_{N_X} and Z_1, \dots, Z_{N_Z} , are in the domain of attraction of the Type I limit
 128 law so that $a_Z(N_Z) = a_X(N_X) = 1$. The following calculation can be straightforwardly
 129 adapted for the eight other possible combinations. Let $G_Z(N_Z) \triangleq Z_{\max} - b_Z(N_Z)$ and
 130 $G_X(N_X) \triangleq X_{\max} - b_X(N_X)$ with density and distribution functions f_{G_Z} , f_{G_X} , F_{G_Z} and
 131 F_{G_X} . We will, for notational convenience, drop the arguments N_Z and N_X . On writing G
 132 for a random variable with the standard Type I distribution, having density and distribution
 133 functions f_G and F_G , we have

$$\begin{aligned}
 p = \text{pr}(Z_{\max} \leq X_{\max}) &= \int_{-\infty}^{\infty} \text{pr}(G_Z \leq v + b_X - b_Z) f_{G_X}(v) dv \\
 &= \int_{-\infty}^{\infty} \{\text{pr}(G_Z \leq v + b_X - b_Z) - \text{pr}(G \leq v + b_X - b_Z)\} f_{G_X}(v) dv \\
 &+ \int_{-\infty}^{\infty} \text{pr}(G \leq v + b_X - b_Z) \{f_{G_X}(v) - f_G(v)\} dv \\
 &+ \int_{-\infty}^{\infty} \text{pr}(G \leq v + b_X - b_Z) f_G(v) dv \triangleq I_1 + I_2 + I_3. \tag{5}
 \end{aligned}$$

134 This illustrates that if the extreme value limit laws are used to approximate the probability
 135 $p = \text{pr}(Z_{\max} \leq X_{\max})$, then the error incurred is given by the sum of the integrals $I_1 + I_2$.
 136 The relevant form of convergence for this type of problem is therefore an appropriately
 137 weighted $\mathbb{L}_1(\text{Leb})$ norm of the density and distribution functions. This is stronger than
 138 uniform convergence of distribution functions, which has been studied for several starting
 139 distributions F_X . Notably, Hall (1979) showed that for the maxima of standard normally
 140 distributed random variables, the uniform convergence rate in (4) to the Type I limit is no
 141 better than $(\log \log n)^2 / \log n$.

142 The conclusion is that, while these limiting approximations are appealing in that they
 143 deliver simple easily interpretable solutions, their adequacy in the present context, partic-
 144 ularly for small N_X or N_Z is not guaranteed and depends heavily on the the distributions
 145 of X and Z . The proposal discussed in §2 provides a compromise between simplicity and

146 adequacy of the resulting analytic approximation.

147 4 An idealized case

148 For an example in which the density function of $V_i = 1 - F_Z(X_i)$ is exactly of the form
 149 $(1 - \gamma)v^{-\gamma}$ used above, let X_i be exponentially distributed of rate ξ and Z_i be exponentially
 150 distributed of rate λ . Then $F_Z^{-1}(z) = -\log(1 - z)/\lambda$ so that

$$\text{pr}(V_i \leq v) = 1 - F_X\{F_Z^{-1}(1 - v)\} = v^{\xi/\lambda}, \quad 0 < v < 1.$$

151 It follows that the density function of each V_i is given by $(1 - \gamma)v^{-\gamma}$, where $\gamma = 1 - \xi/\lambda$.
 152 Thus the probability p is approximated by formulae 1 or 2, or a truncated version thereof,
 153 with $\gamma = 1 - \xi/\lambda$.

154 A direct calculation for the exponentials would entail solving the integral

$$p = N_X \xi \int_0^\infty \{1 - \exp(-\lambda m)\}^{N_Z} \{1 - \exp(-\xi m)\}^{N_X - 1} dm,$$

155 which does not appear to have an exact analytic solution. An alternative is to approximate
 156 p using the Fisher and Tippett Type I limiting form of the rescaled exponential maxima.
 157 This is term I_3 in equation (5) and thus incurs the error $I_1 + I_2$. The scaling constants are

$$b_Z = b_Z(N_Z) = F_Z^{-1}(1 - N_Z^{-1}) = \frac{-\log\{1 - (1 - N_Z^{-1})\}}{\lambda} = \frac{\log N_Z}{\lambda}$$

158 and similarly for b_X . Thus, on letting $B = \xi^{-1} \log(N_X) - \lambda^{-1} \log(N_Z)$,

$$I_3 = \int_{-\infty}^\infty \exp\{-e^{-(v+B)}\} \exp\{-(v + e^{-v})\} dv = \frac{e^B}{1 + e^B} = \frac{N_X^{1/\xi}}{N_Z^{1/\lambda} + N_X^{1/\xi}}. \quad (6)$$

159 The simulations in §5.1 show that this approximation is considerably less accurate than
 160 formula (1) even for very large values of N_X and N_Z , suggesting that at least one of the two
 161 error terms I_1 and I_2 from equation (5) decays slowly for exponentially distributed random
 162 variables. The simpler formula (3) fails because the representation $(1 - \gamma)v^{-\gamma}$ does not
 163 hold even as an approximation for the density function of $W_i \triangleq F_X(Z_i)$ in this example.

5 Numerical assessment

5.1 Empirical analysis of an idealized case

In each of 10000 Monte Carlo replications we generated $N_Z = 100$ random variables Z_1, \dots, Z_{N_Z} from an exponential distribution of rate $\lambda \in \{1.1, 1.2, \dots, 2\}$ and $N_X = 20$ random variables X_1, \dots, X_{N_X} from an exponential distribution of rate $\xi = 1$. The maxima of these two sets of random variables, Z_{\max} and X_{\max} were recorded. The simulated probability that X_{\max} exceeds Z_{\max} was obtained by averaging the indicator random variables $\mathbb{I}(X_{\max} > Z_{\max})$ over the Monte Carlo replications. The maximum likelihood estimate $\hat{\gamma}$ was also obtained in each simulation by fitting a density of the form $(1-\gamma)v^{-\gamma}$ to $V_i = 1 - F_Z(X_i)$ for $i = 1, \dots, N_X$. These maximum likelihood estimates were then averaged over Monte Carlo replications and used in the formula (1). These are plotted against λ in the left panel of Figure 1 along with the version that uses the exact value $\gamma = 1 - \xi/\lambda$. The latter approximation is essentially exact.

The experiment was repeated for $N_X = 50$ and the results are shown in the right panel of Figure 1. As expected, the quality of the analytic approximation is unaffected but the quality of the approximation based on the maximum likelihood estimate of γ is improved because the bias in the maximum likelihood estimates decreases with increasing N_X .

We also report the approximation (6) based on on Fisher's and Tippett's (1928) limit laws. These are evidently inaccurate for the sample sizes under consideration, even though the true values of ξ and λ are used. Further simulations (not reported) indicate that the approximation (6) becomes more accurate as N_X and N_Z increase, but remains rather poor even when $N_X = N_Z = 10^5$.

Figure 2 illustrates the sensitivity of approximation (1) to N_Z for two different values of N_X . Thus, although the formula uses the notional limiting operation $N_Z \rightarrow \infty$ in Stirling's formula, the approximation (1) is accurate even for small N_Z . This is expected in view of the remarkable accuracy of Stirling's approximation to $\Gamma(x)$ even for small x .

5.2 Violation of the main assumption

If the model $(1-\gamma)v^{-\gamma}$ for the density function of the random variables $V_i = 1 - F_Z(X_i)$ is not satisfied to an adequate order of approximation, formula (1) is likely to give inaccurate

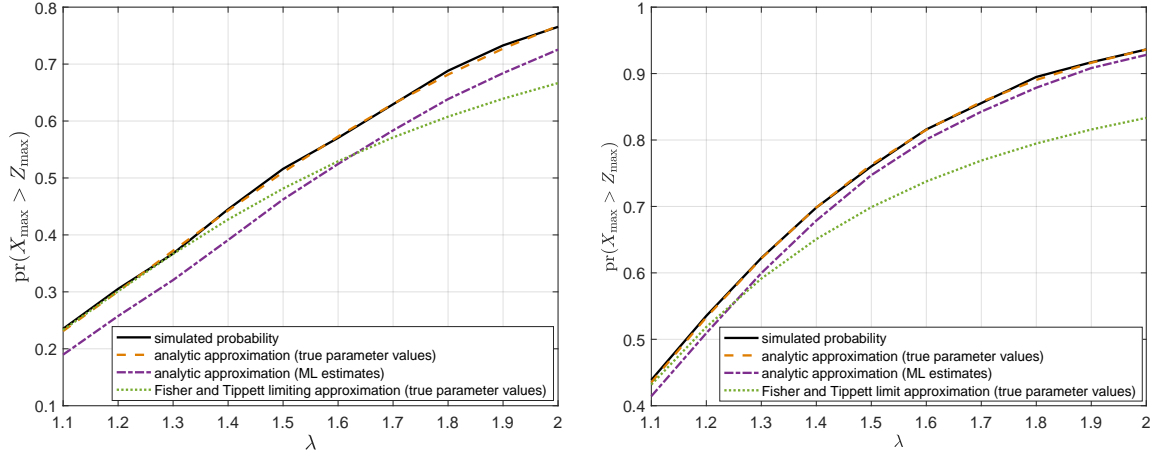


Figure 1: Simulated exceedance probabilities and the near-exact formula based on equation (1) for $N_Z = 100$, $\xi = 1$ and different values of λ , and for $N_X = 20$ (left) and $N_X = 50$ (right). Also depicted is formula (1) with ξ and λ replaced by the average of their maximum likelihood estimates over the 10^5 Monte Carlo replications and the limiting approximation (6) using the true values of ξ and λ .

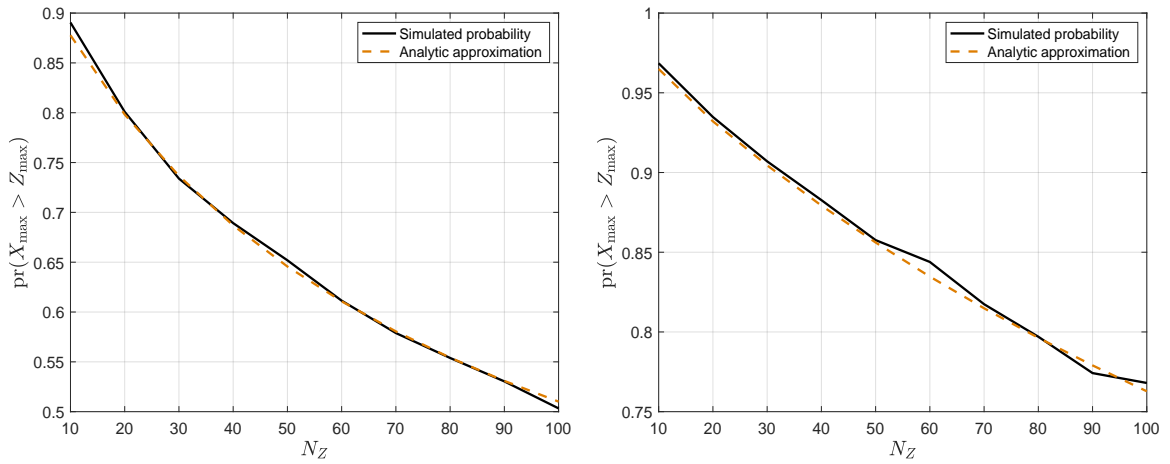


Figure 2: Simulated exceedance probabilities and the exact formula based on equation (1) for $\xi = 1$, $\lambda = 1.5$ and different values of N_Z , and for $N_X = 20$ (left) and $N_X = 50$ (right).

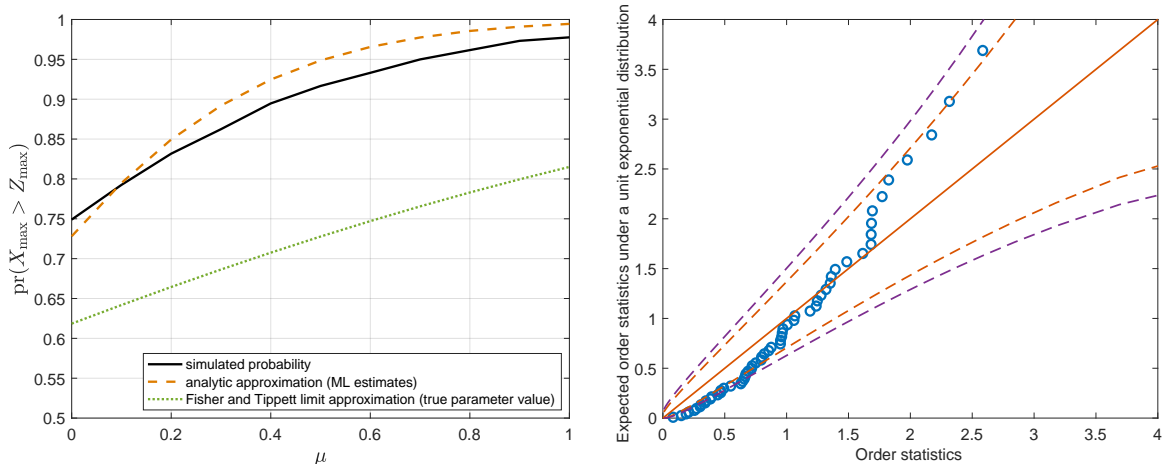


Figure 3: Left: simulated exceedance probabilities and the formula based on equation (1) for an example in which the key assumption is violated, for $N_Z = 20$ and $N_X = 60$. Also displayed is the approximation based on Fisher’s and Tippett’s limit laws. Right: probability plot of $-\log V_i^{1-\hat{\gamma}}$ against the unit exponential order statistics.

193 conclusions. The quality of the approximation can be assessed statistically for large N_X at
 194 a suitable confidence level α by fitting the density $(1 - \gamma)v^{-\gamma}$ by maximum likelihood and
 195 comparing the realization of $2(1 - \hat{\gamma}) \sum_{i=1}^n \log V_i$ to the α upper quantile of a χ^2 distribution
 196 with $2N_X$ degrees of freedom. For an informal indication the order statistics of $-(1 - \hat{\gamma}) \log V_i$
 197 can be plotted against the unit exponential order statistics. This is illustrated in Figure 3.

198 The experiment is as described in the previous section except that Z_1, \dots, Z_{N_Z} are
 199 standard normally distributed and X_1, \dots, X_{N_X} are normally distributed of unit variance
 200 and means μ as displayed on the axis of Figure 3 (left). The values of N_Z and N_X are
 201 20 and 60. The parameter γ of the representation $(1 - \gamma)v^{-\gamma}$ is estimated by maximum
 202 likelihood and used in the formula (1). Also depicted is the limiting approximation I_3
 203 from equation (5). As in equation (6) this approximation is $e^B(1 + e^B)^{-1}$ but with $B =$
 204 $\mu + \Phi^{-1}(1 - N_X^{-1}) - \Phi^{-1}(1 - N_Z^{-1})$, where Φ is the standard normal distribution function. The
 205 true value of μ is used in this latter approximation, yet the approximation I_3 is poor. This
 206 is unsurprising because, even for the weaker uniform convergence of distribution functions,
 207 convergence rates are extremely slow for normal distributions, as discussed in §3. The
 208 approximation (1) is less inaccurate and qualitatively successful, in spite of considerable
 209 violation of the assumption used in the calculation, as indicated by the right hand panel of
 210 Figure 3. This is a plot of the order statistics of $-(1 - \hat{\gamma}) \log V_i$ against the unit exponential

211 order statistics.

212 6 Conclusion

213 We have illustrated how the transformation to $U_i = 1 - F_Z(Z_i)$ and $V_i = 1 - F_Z(X_i)$
214 facilitates probability calculations involving order statistics for the two populations. The
215 accuracy of the ensuing approximations hinges of the plausibility of the $(1 - \gamma)v^{-\gamma}$ assump-
216 tion for the probability density function of each V_i . This can be assessed as in §5.2.

217 We have not discussed statistical aspects associated with the various complicating sce-
218 narios that could be envisaged, for instance if F_Z needs to be estimated. If there was an
219 auxiliary sample in which the class labels were known, the simplest nonparametric estimator
220 is $\hat{F}_Z(z) = n^{-1} \sum_{i=1}^{N_Z} \mathbb{I}\{X_i \leq z\}$, leading to

$$\hat{U}_i \triangleq 1 - \hat{F}_Z(Z_i) = U_i + F_Z(Z_i) - \hat{F}_Z(Z_i).$$

221 The final term is bounded in absolute value by $\sup_{z \in \mathbb{R}} |\hat{F}_Z(z) - F_Z(z)|$. We have

$$\text{pr}(\sup_{z \in \mathbb{R}} |\hat{F}_Z(z) - F_Z(z)| < \varepsilon) \geq 1 - 2e^{-2N_Z\varepsilon^2}$$

222 (Dvoretzky, Kiefer and Wolfowitz, 1956; Massart, 1990) implying (see §3.8 of Barndorff-
223 Nielsen and Cox, 1989) that $\hat{U}_i = U_i + O_p(N_Z^{-1/2})$. The same argument applies for the
224 V_1, \dots, V_{N_X} components.

225 **Acknowledgement:** the work was supported by the EPSRC (EP/P002757/1). I am
226 grateful to the referees for helpful comments and Nicholas Beale of Sciteb Ltd. for asking a
227 question that led to some of these calculations.

228 REFERENCES

229 Barndorff-Nielsen, O. E. and Cox, D. R. (1989). *Asymptotic Techniques for Use in*
230 *Statistics*. Chapman and Hall, London.

231

232 Bairamov, I. and Parsi, S. (2011). Order statistics from mixed exchangeable random

233 variables. *J. Comput. Appl. Math.*, 235, 4629–4638.

234

235 Bayramoglu (Bairamov), I. and Eryilmaz, S. (2015). Order statistics of dependent se-
236 quences consisting of two sets of exchangeable variables. *J. Comput. Appl. Math.*, 286, 1–6.

237

238 Dvoretzky, A. Kiefer, J. and Wolfowitz, J. (1956). Asymptotic minimax character of the
239 sample distribution function and of the classical multinomial estimator. *Ann. Math. Statist.*,
240 27, 642–669.

241

242 Fisher, R. A. and Tippett, L. H. C. (1928). Limiting forms of the frequency distribution
243 of the largest or smallest of a sample. *Proc. Camb. Phil. Soc.*, 24, 180–90.

244

245 Gnedenko, B. V. (1943). Sur la distribution limite du terme maximum d’une série aléa-
246 toire, *Ann. Math.*, 44, 423–53.

247

248 Gumbel, E. J. and von Schelling, H. (1950). The distribution of the number of ex-
249 ceedances. *Ann. Math. Statist.*, 21, 247–262.

250

251 Hall, P. G. (1979). On the rate of convergence of normal extremes. *J. App. Probab.*, 16,
252 433–39.

253

254 Massart, P. (1990). The tight constant in the Dvoretzky–Kiefer–Wolfowitz inequality.
255 *Ann. Probab.*, 18, 1269–1283.

256

257 Sarkadi, K. (1957). On the distribution of the number of exceedances. *Ann. Math. Statist.*,
258 28, 1021–1023.

259